

## روش‌های تحلیل رگرسیونی برای داده‌های بعد بالا

منیره معنوی<sup>۱</sup>، مهدی روزبه<sup>۲</sup>

تاریخ دریافت: ۹۹/۴/۳۱

تاریخ پذیرش: ۹۹/۱۱/۱

چکیده:

با پیشرفت علم، دانش و تکنولوژی، روش‌های جدید و جامع برای اندازه‌گیری، جمع‌آوری و ثبت اطلاعات ابداع شده‌اند، که منجر به ظهور و توسعه داده‌های بعد بالا شده‌اند. مجموعه داده‌های بعد بالا، یعنی مجموعه داده‌هایی که در آن تعداد متغیرهای توضیحی بسیار بزرگ‌تر از تعداد مشاهدات است، به سادگی و با روش‌های سنتی و کلاسیک، مانند روش کمترین توان‌های دوم معمولی، نمی‌توانند تحلیل شوند و تفسیرپذیری آن امری بسیار پیچیده خواهد بود. اگرچه در صورتی که فرضیات اساسی برقرار باشند، برآورد کمترین توان‌های دوم معمولی بهترین روش برآورد در تحلیل رگرسیونی است ولی برای داده‌های بعد بالا قابل استفاده نبوده و در این شرایط مستلزم به کارگیری روش‌هایی نوینی هستیم. در این مقاله در ابتدا، به مشکلات روش‌های کلاسیک در تحلیل داده‌های بعد بالا اشاره می‌شود و سپس، به معرفی و توضیح روش‌های تحلیل رگرسیونی متداول و امروزی مانند روش‌های تحلیل مؤلفه اصلی و توانیده برای داده‌های بعد بالا پرداخته می‌شود. در انتها یک مطالعه شبیه‌سازی و تحلیل داده واقعی برای بررسی و مقایسه روش‌های اشاره شده در داده‌های بعد بالا انجام می‌گردد.

**واژه‌های کلیدی:** تحلیل مؤلفه‌های اصلی، مجموعه داده‌های بعد بالا، روش کمترین توان‌های دوم توانیده.

با  $d$  رشد کند، تا بتواند قابل تحلیل و تفسیر با روش‌های کلاسیک باشد.

### ۱ مقدمه

در گذشته، کار با داده‌های کلاسیک که دارای حداکثر چند ده متغیر توضیحی بوده‌اند، بسیار ساده بود و روش‌های کلاسیک نظیر کمترین توان‌های دوم نتایج قابل قبولی ارائه می‌دادند. اما این روش‌ها در حضور داده‌های بعد بالا برآوردی ارائه نمی‌دهند. در تحلیل رگرسیونی داده‌های بعد بالا تعداد زیاد متغیرهای توضیحی، محقق را با چالشی جدی روبه‌رو می‌کند و همیشه جدالی بین دقت، سرعت و هزینه وجود دارد. از سوی دیگر، وجود متغیرهایی در مدل که با متغیر پاسخ ارتباطی ندارند، منجر به بیش‌برازشی<sup>۳</sup> می‌شود. چنین مدلی، مدلی بسیار پیچیده برای داده‌ها است. در این شرایط، مدل با تغییرات جهشی سعی در پوشش داده‌های حاصل از نمونه و حتی مقدارهای خطاها می‌کند. درحالی‌که این مدل باید منعکس‌کننده رفتار جامعه باشد. به کارگیری تمامی متغیرهای توضیحی زمانی طولانی صرف کرده و هزینه‌های محاسباتی بسیاری را بر محقق تحمیل می‌کند. در این شرایط این ابهام مطرح می‌شود که شاید حضور تمامی متغیرها برای برازش مدل الزامی نباشد و بدین‌سان اغلب محققان با در نظر گرفتن فرض تنگی<sup>۵</sup> سعی در کاهش بعد دارند و از زیرمجموعه‌ای مناسب از متغیرها برای برازش

امروزه با گسترش روزافزون علم، دانش و فناوری، روش‌های نوین و دقیقی برای اندازه‌گیری، جمع‌آوری و ثبت اطلاعات ابداع شده و این امر باعث ظهور و گسترش داده‌های بعد بالا شده است. از اواخر سال ۱۹۹۰ کار با این مجموعه داده‌ها یعنی داده‌هایی که در آن تعداد متغیرهای توضیحی ( $p$ ) بسیار بیشتر از تعداد مشاهدات ( $n$ ) است، شروع شد [۱۰]. افزایش بیش‌ازپیش ابعاد داده‌ها به یک مسئله اساسی در مبحث داده‌کاوی تبدیل شده است. تحلیل و تفسیر این داده‌ها امری دشوار و بسیار قابل‌تأمل است به طوری که واسرمن [۲۵] در کتاب خود به این امر اشاره نموده و از این پدیده تحت عنوان نفرین ابعاد<sup>۳</sup> یاد کرده است. اصطلاحی که معمولاً به بلمن [۶] نسبت داده می‌شود. این بدان معنی است که با افزایش بعد متغیرها، برآورد پارامترهای مربوطه بسیار دشوار می‌شود. حداقل دو نوع از این نفرین وجود دارد. اولین مورد، نفرین محاسباتی ابعاد است. این امر به این واقعیت اشاره دارد که بار محاسباتی برخی از روش‌ها می‌تواند به صورت نمایی با بعد افزایش یابد. دومین مورد که واسرمن آن را نفرین آماری ابعاد نامید به این صورت است که اگر داده‌ها دارای بعد  $d$  باشند، پس نمونه‌ای با اندازه  $n$  موردنیاز است که به صورت نمایی

<sup>۱</sup>دانش‌آموخته کارشناسی ارشد آمار، دانشگاه سمنان، سمنان، ایران.

<sup>۲</sup> هیئت‌علمی گروه آمار، دانشگاه سمنان، سمنان، ایران mahdi.roozbeh@semnan.ac.ir

<sup>۳</sup> Curse of Dimensionality

<sup>۴</sup> Overfitting

<sup>۵</sup> Sparsity

زیرمجموعه، کوچک‌ترین زیرمجموعه‌ای از متغیرها است که به اندازه کافی تغییرات متغیر پاسخ را توجیه می‌کند. برای تعیین مدل با این روش ابتدا مدلی که تنها حاوی عرض از مبدأ بوده در نظر گرفته می‌شود. سپس تمامی زیرمجموعه‌های تک عضوی و دو عضوی، سه عضوی و ... در نظر گرفته می‌شود. و پس از آن مدل‌های رگرسیونی ایجاد شده و برازش داده می‌شوند. در قدم آخر تمامی مدل‌های ایجاد شده با یکی از معیارهای سنجش نیکویی برازش مدل نظیر ضریب تعیین، تابع زیان، مجموع توان‌های دوم خطا،  $C_p$  مالوس<sup>۸</sup>، اعتبارسنجی متقابل<sup>۹</sup> و ... مقایسه می‌شوند و بهترین آن‌ها به عنوان مدل نهایی انتخاب می‌شود.

در این قسمت سؤال بسیار مهمی پیش می‌آید: اندازه زیرمجموعه‌ها در چه بازه‌ای تغییر می‌کند؟ این سؤال به سادگی از طریق رابطه

$$\min\{0, 1, \dots, M\}, \quad M \leq \min\{n-1, (p+1)\}$$

پاسخ داده می‌شود. که در آن،  $M$  اندازه زیرمجموعه‌ها،  $n$  تعداد مشاهدات و  $p$  تعداد متغیرهای توضیحی است.

رگرسیون زیرمجموعه ضرایب متغیرهایی که در زیرمجموعه انتخابی وجود نداشته باشند، صفر در نظر می‌گیرد و یا به عبارت دیگر آن را متورم می‌کند [۹]. این روش، به دو دلیل کاهش واریانس و سادگی تفسیر نتایج (برای تعداد متغیرهای توضیحی کم) مورد توجه قرار گرفته است [۹]. رقیب اول رگرسیون زیرمجموعه از نظر کاهش واریانس، روش ستیغی [۱۷] است. رگرسیون زیرمجموعه می‌تواند دقت پیش‌گویی را افزایش دهد اما اگر تعداد متغیرهای موجود در مدل،  $M$ ، بزرگ باشد، این مسئله ایهام‌برانگیز است.

یک جنبه منفی رگرسیون بهترین زیرمجموعه بی‌ثباتی آن در اثر انحراف‌های کوچک در داده‌هاست. به عنوان مثال، اگر تنها مشاهده  $(y_n, x_n)$  از مجموعه داده‌ها حذف شوند و از روش‌های انتخاب یکسان استفاده شود دو زیرمجموعه نهایی که با استفاده از این مشاهده و بدون استفاده از آن تعیین شده‌اند، بسیار متفاوت خواهند بود. بنابراین منجر به تغییر شدید در معادله پیش‌گویی خواهد شد. این در حالی است که اگر از روش ستیغی استفاده شود (البته با پارامتر ستیغی یکسان) برآوردگرهای حاصل از حذف مشاهده  $(y_n, x_n)$  و برآوردگر حاصل بدون حذف آن تفاوت چندانی نخواهند داشت [۹].

یکی دیگر از نقاط ضعف این روش زمانی است که تعداد متغیرهای توضیحی بیشتر از ۴۰ باشد. در این شرایط به دلیل پیچیدگی ترکیباتی، رگرسیون بهترین زیرمجموعه به یک مسئله دشوار و پیچیده تبدیل می‌شود،

مدل کمک می‌گیرند و از این رو علاقه‌مند به شناسایی این زیرمجموعه مناسب هستند [۷]. حتی امروزه مدل‌های رگرسیونی لوژستیک نیز در این زمینه مورد اهمیت ویژه‌ای قرار گرفته‌اند [۲۴].

**مثال ۱.۱.** در مطالعات ژنتیکی که در اکثر آن‌ها ارتباط بین تغییرات  $DNA$  و بیماری هدف مورد مطالعه می‌باشد، ۹۹٫۸ درصد از سکانس  $DNA$  انسان‌ها با یکدیگر یکسان است و بیش از ۸۰ درصد از آن ۰٫۱ درصد باقی‌مانده به  $SNP$  (پلی مورفیسم تک نوکلئوتیدی)<sup>۶</sup> مرتبط است. میلیون‌ها  $SNP$  با هم تعامل دارند تا فنوتیپ نهایی بیماری را تعیین کنند. این تعداد زیاد  $SNP$  ها و امکان وجود اثرهای متقابل بین عوامل ژنتیکی با هم و با عوامل محیطی یک چالش را در تحلیل این که کدام یک از این  $SNP$  ها در ارتباط با بیماری هستند، پدید آورده است. هدف اصلی در تحلیل این اطلاعات تعیین این که کدام  $SNP$  یا مجموعه‌ای از آن‌ها روی بیماری مؤثر هستند، نیست چرا که بسیاری از  $SNP$  ها اثرهای حاشیه‌ای ناچیزی دارند ولی اثرات متقابل آن‌ها بسیار قوی است. هدف اصلی اولویت‌بندی  $SNP$  ها بر اساس اهمیت و نقش آن‌ها در ارتباط با بیماری به منظور مطالعه و بررسی بیشتر آن‌هاست. در این شرایط برای رسیدن به این هدف در نظر گرفتن فرض تنگی الزامی است [۴].

در ادامه به معرفی مدل رگرسیونی پیردازیم. مدل رگرسیونی خطی ساده به صورت

$$y = X\beta + \varepsilon, \quad (1)$$

است، به طوری که در آن  $y = (y_1, \dots, y_n)^T$  بردار متغیر پاسخ،  $X = (x_1, \dots, x_p)_{n \times p}$  ماتریس مشاهدات متغیرهای توضیحی با  $x_i = (x_{1i}, \dots, x_{ni})^T$  بردار پارامترها و  $\beta = (\beta_1, \dots, \beta_p)^T$  بردار پارامترها و  $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)^T$  بردار خطا است.

## ۲ روش‌های تحلیل رگرسیونی برای داده‌های بعد بالا

در این بخش، چندین روش برای رویارویی با داده‌های بعد بالا معرفی می‌شود. برخی از این روش‌ها فرض تنگی را در نظر می‌گیرند و برخی در برنمی‌گیرند.

### ۱.۲ روش‌های کلاسیک

رگرسیون بهترین زیرمجموعه<sup>۷</sup> یکی از روش‌های کلاسیک است. ایده اصلی این روش، ساخت مدلی با زیرمجموعه‌ای از متغیرها است. در اصل این

<sup>۶</sup>Single-nucleotide polymorphism

<sup>۷</sup>Best-subset regression

<sup>۸</sup> $C_p$ -Mallows

<sup>۹</sup>Cross validation

مدل بسیار متفاوت خواهد شد و دقت پیش‌گویی کاهش می‌دهند. بنابراین باید در صحت پیش‌گویی‌ها تردید نمود. برای برطرف کردن این مشکلات، می‌توان از روش‌های دیگری که در ادامه مطرح می‌شود، کمک گرفت [۲].

## ۲.۲ روش مؤلفه‌های اصلی

روش مؤلفه‌های اصلی<sup>۱۴</sup>، روشی بسیار ساده و کاربردی است که برای نخستین بار توسط پیرسون [۲۰] ریاضی‌دان انگلیسی مطرح شد و در سال ۱۹۳۳ به‌طور مستقل توسط هتلینگک به‌منظور تحلیل ساختارهای ماتریس‌های واریانس کوواریانس و ضریب همبستگی توسعه داده شد [۱۸]. ایده اصلی این روش کاهش ابعاد داده‌ها به نحوی است که حداکثر اطلاعات ممکن موجود در داده‌ها محفوظ بماند. همانند بسیاری از روش‌های چندمتغیره تا قبل از ظهور رایانه و نرم‌افزارهای مرتبط به دلیل پیچیدگی محاسبات به‌صورت گسترده مورداستفاده واقع نشد. اما اکنون تقریباً در تمامی نرم‌افزارهای آماری و ریاضیاتی بسته‌های محاسباتی آن توسط محققان تدوین شده است. طبیعتاً پس از پیدایش نرم‌افزارها و بسته‌های رایانه‌ای این روش بسیار محبوب و رایج گردید.

هدف اصلی این روش یافتن ترکیب‌های خطی از متغیرهای توضیحی است که واریانس را به بیشینه خود برساند. این روش بر پایه مقادیر ویژه و بردارهای ویژه ماتریس واریانس-کوواریانس یا ماتریس همبستگی استوار است. روش مؤلفه اصلی، معمولاً یک تجزیه نهایی تلقی نمی‌شود بلکه به‌عنوان ابزاری میانی برای مطالعه و بررسی‌های بیشتر مورداستفاده قرار می‌گیرد.

از دید ریاضیات این روش، یک تبدیل خطی متعامد است که داده‌ها را به دستگاه مختصات جدید می‌برد به‌طوری‌که بزرگ‌ترین واریانس داده‌ها بر روی اولین محور مختصات، دومین بزرگ‌ترین واریانس بر روی دومین محور مختصات قرار می‌گیرد و به همین ترتیب برای سایر متغیرها ادامه می‌یابد. این امر از طریق تبدیل متغیرها به متغیرهای جدیدی رخ می‌دهد که مؤلفه‌های اصلی نامیده می‌شوند. مؤلفه‌های اصلی ناهمبسته بوده و به ترتیبی اولویت‌بندی می‌شوند که تعداد اندکی از آن‌ها اغلب تغییرات موجود در متغیرهای اولیه را با خود به همراه دارند.

از کاربردهای دیگر تحلیل مؤلفه‌های اصلی می‌توان به تبدیل متغیرهای همبسته به متغیرهای ناهمبسته، یافتن ترکیبات خطی با تغییرپذیری نسبی بزرگ یا کوچک، کاهش حجم داده‌ها و تفسیر ساده‌تر آن‌ها اشاره نمود.

گام بعدی پس از ساخت مؤلفه‌ها تعیین تعداد مؤلفه‌هایی است که برای

زیرا  $2^p$  زیرمجموعه وجود دارد که حل آن هزینه‌های محاسباتی زیادی را به همراه دارد. حتی اگر محاسبات نیز انجام‌پذیر باشد، با چنین فضای جست‌وجوی عظیمی، واریانس این روش بسیار بالا بوده و این امر بدین معناست که این روش تنها برای  $p$  های نسبتاً کوچک کاربردی است.

روش دیگری که در این بخش مطرح می‌شود، رگرسیون گام‌به‌گام<sup>۱۰</sup> نام دارد. رگرسیون گام‌به‌گام به سه نوع پس‌رو<sup>۱۱</sup>، پیش‌رو<sup>۱۲</sup> و مرحله‌ای<sup>۱۳</sup> (ترکیبی از هر دو روش پیش‌رو و پس‌رو) تقسیم می‌شود.

روش پس‌رو، مدلی را که در آن تمامی متغیرها حضور دارند در نظر گرفته و با ملاک قرار دادن یکی از معیارهای سنجش اقدام به بررسی برای حذف متغیرهای بی‌تأثیر در مدل می‌نماید.

روش پیش‌رو، مدلی را که تنها در آن عرض از مبدأ وجود دارد در نظر گرفته و با ملاک قرار دادن یکی از معیارهای سنجش اقدام به یافتن تأثیرگذارترین متغیر می‌نماید و آن را به مدل اضافه می‌کند و این روند را تا زمانی که معیار سنجش معنی‌دار نشود، ادامه می‌دهد.

می‌توان گفت که روش مرحله‌ای ترکیبی از روش‌های پس‌رو و پیش‌رو می‌باشد که در آن، در هر مرحله هم‌زمان با حذف یک متغیر از مدل بر اساس معیار سنجش، وجود همان متغیر در مدل بر اساس معیار سنجش بررسی می‌شود [۱۵].

روش رگرسیون گام‌به‌گام پیش‌رو از یک اصلاح ساده در روش رگرسیون بهترین زیرمجموعه به‌دست‌آمده است. در رگرسیون گام‌به‌گام پیش‌رو برخلاف روش رگرسیون بهترین زیرمجموعه، زیرمجموعه‌ها آشیانه‌ای (تودرتو) هستند و این باعث می‌شود تا روش گام‌به‌گام پیش‌رو قابل کنترل‌تر باشد و برای مدل‌هایی که تعداد متغیرها بیشتر از ۴۰ هست، نیز انجام‌پذیر باشد. اما اگر تعداد متغیرها بسیار زیاد باشد، برای مثال بیش از هزار متغیر توضیحی در مدل حضور داشته باشد، این روش بسیار پیچیده، خسته‌کننده و زمان‌گیر است [۱۰].

امروزه هزینه‌های محاسباتی بسیار زیاد روش‌های کلاسیک، تحلیل رگرسیونی داده‌های بعد بالا را با چالشی بسیار جدی مواجه ساخته است و شیرینی سادگی تفسیرپذیری مدل را به کام محققان تلخ می‌کند. این روش‌ها از منظر محاسباتی اصلاً مقرون‌به‌صرفه نیستند و همان‌طور که ذکر شد، محاسبات این روش‌ها ذکر شده عموماً بسیار پیچیده، گیج‌کننده و زمان‌بر است. اما مهم‌ترین اشکال این روش‌ها، ناپایداری آن‌ها است زیرا آن‌ها فرآیندی گسسته هستند که با کوچک‌ترین تغییری در داده‌ها نتیجه انتخاب

<sup>10</sup>Stepwise regression

<sup>11</sup>Backward

<sup>12</sup>Forward

<sup>13</sup>Stagewise

<sup>14</sup>Principle component

توضیحی اولیه ممکن است مقیاس‌های گوناگونی داشته باشند. به‌عنوان مثال می‌توان به یک مجموعه داده که شامل متغیرهایی با یکاهای گالون، کیلومتر، سال نوری و ... است، اشاره کرد. بدیهی است که مقدار واریانس این متغیرها اعداد بزرگی خواهد بود. اعمال مؤلفه‌های اصلی روی متغیرهای استاندارد نشده منجر به تأثیر عظیمی بر روی متغیرهایی که واریانس بزرگی دارند، می‌شود و این امر به‌نوبه خود می‌تواند منجر به وابستگی مؤلفه اصلی به متغیرهای دارای واریانس بالا شود که بسیار نامطلوب است. البته برای حل این مشکل می‌توان از ماتریس همبستگی نیز کمک گرفت زیرا ماتریس همبستگی برخلاف ماتریس واریانس-کواریانس مقیاس پایاست.

### ۳.۲ روش کمترین توان‌های دوم تاوانیده

برآوردگرهای تاوانیده از به‌کارگیری روش‌های کمترین توان‌های دوم یا ماکسیمم درست‌نمایی تاوانیده در مدل‌بندی رگرسیونی حاصل می‌شود. در حقیقت این برآوردگرها از روش بهینه‌سازی یک تابع درجه دوم (البته نه همیشه) نسبت به یک برآوردگر تاوانیده به دست می‌آید. ایده اصلی در این برآوردگرها، این است که متغیرهای توضیحی کم‌اهمیت در یک مدل تنگ به سمت صفر منقبض شود. در اکثر روش‌های تاوانیده این متغیرها از مدل خارج می‌شوند و همین امر موجب ساده‌تر شدن تفسیرپذیری مدل خواهد شد [۱].

تابع هدف در روش کمترین توان‌های دوم معمولی در مدل رگرسیونی چندگانه (۱) به صورت زیر است:

$$\min_{\beta \in \mathbb{R}^p} \left\{ (y - X\beta)^T (y - X\beta) \right\}. \quad (2)$$

با محاسبه مشتق تابع هدف عبارت (۲) نسبت به بردار پارامتر  $\beta$  و برابر صفر قرار دادن آن و حل معادله حاصل، برآوردگر کمترین توان‌های دوم به‌سادگی یافت می‌شود. به عبارت دیگر در این روش خطا برحسب  $\beta$  کمینه می‌شود. واضح است که اگر مقادیر واقعی  $\beta$  بزرگ باشد، برآوردگر حاصل، فارغ از احتساب بزرگی  $\beta$ ، فاصله زیادی از مقدار واقعی  $\beta$  خواهد داشت و به همین سبب خطای برآورد زیاد شده و دقت آن کاهش می‌یابد [۱]. یک روش بسیار سودمند و تقریباً ساده برای مقابله با این مشکل، تاوان دادن مقادیر بزرگ  $\beta$  است بدین سان به جای کمینه کردن تابع هدف روش کمترین توان‌های دوم تابع هدف کمترین توان‌های دوم تاوانیده یعنی

$$\min_{\beta \in \mathbb{R}^p} \left\{ (y - X\beta)^T (y - X\beta) \right\}, \quad s.t. \quad p(\beta) \leq t \quad (3)$$

کمینه می‌شود.  $t \geq 0$  پارامتر تنظیم‌کننده<sup>۱۶</sup> نامیده می‌شود که اگر برابر صفر باشد، مدل تنها شامل عرض از مبدأ است و اگر بی‌نهایت باشد، مدل کامل

ایجاد مدل رگرسیونی مورد استفاده می‌گیرد. در واقع در این مرحله کاهش بعد صورت می‌گیرد. انتخاب تعداد مناسب مؤلفه‌ها امری بسیار ضروری و حائز اهمیت است، چراکه هدف هر محقق در این مرحله نگهداری کمترین تعداد مؤلفه‌ها با حفظ بیشترین اطلاعات و تغییرپذیری موجود در مجموعه داده‌ها به منظور دسترسی به بالاترین دقت با صرف کمترین زمان ممکن است. روش‌های متعددی برای تعیین تعداد مناسب مؤلفه‌ها وجود دارد که در ادامه به برخی از این موارد به اختصار اشاره می‌شود.

یکی از روش‌های تعیین مؤلفه‌ها توجه به تغییرپذیری بیان‌شده توسط مؤلفه‌ها است. در واقع در این روش تعدادی از مؤلفه‌ها که بتوانند درصد قابل توجهی از واریانس (تغییرپذیری) کل را توجیه کنند، انتخاب می‌شوند. البته باید توجه داشت که لازم است این درصد نسبت به تعداد مؤلفه‌ها به صرفه باشد. برای مثال اگر  $100\%$  مؤلفه  $20\%$  درصد تغییرات واریانس کل و با اتخاذ  $40\%$  مؤلفه  $75\%$  درصد تغییرات واریانس کل توجیه شود، افزودن  $20\%$  مؤلفه دیگر برای دست یافتن تنها به  $5\%$  درصد تغییرات بیشتر به صرفه نخواهد بود. ولی اگر با در نظر گرفتن  $23\%$  مؤلفه  $80\%$  درصد تغییرات واریانس کل توجیه شود، انتخاب بسیار مناسب و معقولی خواهد بود.

روش دیگر، انتخاب تعداد مؤلفه‌هایی است که واریانس آن‌ها بزرگ‌تر یا مساوی متوسط واریانس کل است. لازم به ذکر است که اگر از ماتریس همبستگی استفاده شود، مؤلفه‌هایی که واریانس آن‌ها بزرگ‌تر یا مساوی ۱ است، انتخاب می‌شوند.

علاوه بر دو روش ذکر شده روش‌های شهودی جالبی نیز وجود دارند. یکی از این روش‌ها استفاده از نمودار بازو<sup>۱۵</sup> است. در این نمودار مقادیر ویژه هر مؤلفه ( $\lambda_i$ ) در برابر  $i$  رسم می‌شود. با مشخص کردن ناحیه‌ای از نمودار که شیب آن به‌طور ناگهانی کند شده و شکسته می‌شود، تعداد مؤلفه‌ها تعیین می‌شود. البته بعضی محققان لگاریتم  $\lambda_i$  در برابر  $i$  رسم می‌کنند. لگاریتم تابعی یک‌به‌یک است که مقیاس را کاهش می‌دهد بنابراین استفاده از این روش تفاوت چندانی با نمودار بازو ندارد، اما استفاده از نمودار بازو سودمندتر خواهد بود چراکه جهش‌های ناگهانی در شیب نمودار بهتر نمایش داده می‌شود و بدین سان انتخاب تعداد مؤلفه‌ها ساده‌تر خواهد بود.

لازم به ذکر است که ممکن است هر کدام از روش‌های مطرح شده در بالا پاسخ‌های متفاوتی ارائه دهند که بنا به نوع هدف و سلیقه محقق یکی از این روش‌ها برای انتخاب تعداد مناسب متغیرها به کار گرفته می‌شود.

یکی از نکات بسیار ضروری در استفاده از روش مؤلفه‌های اصلی استانداردسازی متغیرهاست. اهمیت این امر بدین سان است که متغیرهای

<sup>15</sup>Scree diagram

<sup>16</sup>Tuning parameter

برازش و پیچیدگی مدل (حضور متغیرهای زیاد) را برقرار می کند [۱۹]. در ادامه به معرفی چندین روش تاوانیده پرداخته می شود.

## ۴.۲ رگرسیون بریج

نخستین بار فرانک و فریدمن [۱۴] رگرسیون بریج<sup>۲۱</sup> را معرفی نمودند. آن‌ها در این روش، تابع هدف را با در نظر گرفتن محدودیت  $\sum_{j=1}^p |\beta_j|^\gamma \leq t$  کمینه و مسئله بهینه‌سازی را به صورت

$$\min_{\beta \in \mathbb{R}^p} \left\{ \sum_{i=1}^n (y_i - \sum_{j=1}^p \mathbf{X}_{ij} \beta_j)^2 + \lambda \sum_{j=1}^p |\beta_j|^\gamma \right\},$$

$$i = 1, \dots, n, \quad j = 1, \dots, p$$

تدوین کردند که در آن  $\gamma$  مقداری نامنفی است. البته فرانک و فریدمن [۱۴] این مسئله بهینه‌سازی را برای تمامی مقادیر ممکن  $\gamma$  به کار نگرفتند. لازم به ذکر است که به ازای  $\gamma > 1$  ضرایب روش بریج معمولاً برای تمامی متغیرها غیر صفر هستند.

خانواده بریج شامل موارد خاص رگرسیون بهترین زیرمجموعه به ازای  $\gamma = 0$ ، رگرسیون لاسو به ازای  $\gamma = 1$  و رگرسیون ستیغی به ازای  $\gamma = 2$  می‌باشد که در شکل ۱ رگرسیون بریج به ازای این مقادیر رسم شده است.

یک ویژگی بسیار سودمندی که برخی از روش‌های تاوانیده دارند، محدب بودن مسئله بهینه‌سازی آن‌هاست. مسئله بهینه‌سازی روش بریج به ازای  $\gamma \geq 1$  محدب و به ازای  $\gamma < 1$  غیر محدب است [۱۹].

برای محاسبه برآوردگر کمترین توان‌های دوم یعنی  $\hat{\beta}_{OLS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$  نیاز به محاسبه وارون ماتریس  $\mathbf{X}^T \mathbf{X}$  است. ولی این ماتریس در مجموعه داده‌های با بعد بالا به سبب  $p \gg n$  رتبه کامل نبوده و در نتیجه وارون‌پذیر نخواهد بود. البته می‌توان برای یافتن وارون ماتریس  $\mathbf{X}^T \mathbf{X}$  از روش وارون تعمیم‌یافته نیز استفاده نمود اما در این روش وارون به دست آمده منحصر به فرد نبوده و به دنبال آن پاسخ‌های به دست آمده نیز منحصر به فرد نخواهد بود. گاهی اوقات وارون وجود دارد ولی کران‌دار نیست و بدین سبب لازم است از روش‌های دیگری استفاده شود. روش ستیغی مقدار مثبتی مانند  $\lambda$  را به درایه‌های قطر اصلی ماتریس  $\mathbf{X}^T \mathbf{X}$  می‌افزاید و با انجام این عمل می‌توان وارون ماتریس را به سادگی یافت. البته این امر موجب تحریف داده‌ها

(با حضور تمامی متغیرها) خواهد بود.  $p(\beta)$  نیز تابعی از بردار پارامتر است که تابع تاوان<sup>۱۷</sup> نامیده شده و با توجه به نوع آن روش‌های تاوانیده متفاوتی ایجاد می‌شود. برای یافتن کمینه (بیشینه) یک تابع چند متغیره که با یک یا چند محدودیت مواجه است، روش لاگرانژ به کار گرفته می‌شود و با استفاده از روش لاگرانژ چند متغیره عبارت (۳) به صورت زیر بازنویسی می‌شود [۱۱]:

$$\min_{\beta \in \mathbb{R}^p} \left\{ (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) + \lambda p(\beta) \right\}. \quad (4)$$

در عبارت (۴) پارامتر نامنفی  $\lambda$ ، پارامتر منظم‌سازی<sup>۱۸</sup> یا پارامتر تاوان<sup>۱۹</sup> نامیده می‌شود. اگر  $\lambda$  بی‌نهایت باشد، تنها عرض از مبدأ در مدل حضور خواهد داشت و اگر برابر صفر باشد، قسمت تاوان از مدل حذف شده و مدل کمترین توان‌های دوم حاصل خواهد شد.

در اصل،  $\lambda$  و  $t$  نقش یکسانی در مدل ایفا می‌کنند و کنترل بزرگی یا کوچکی ناحیه محدودیت (تاوان) بر دوش این دو پارامتر است. جالب توجه است که رابطه این دو پارامتر عکس یکدیگر است. لازم به ذکر است که در دو عبارت (۳) و (۴) از تابع زیان توان دوم خطا<sup>۲۰</sup> استفاده شده است که بیانگر روش کمترین توان‌های دوم تاوانیده است. اما به طور کلی مسئله بهینه‌سازی تاوانیده به صورت

$$\hat{\beta} = \min_{\beta \in \mathbb{R}^p} \left\{ L(\mathbf{y}, \mathbf{X}\beta) + \lambda p(\beta) \right\}$$

نوشته می‌شود که در آن  $L(\mathbf{y}, \mathbf{X}\beta)$  تابع زیان نامنفی برای نیکویی برازش،  $p(\beta)$  تابع تاوان نامنفی و  $\lambda$  پارامتر منظم‌سازی است که تعادل بین نیکویی

## ۵.۲ رگرسیون ستیغی

آندری نیکولایویچ تیخونوف<sup>۲۲</sup> متخصص ژئوفیزیک و ریاضی‌دان روسی که در زمینه کاربردهای ریاضی در فیزیک کار می‌کرد، منظم‌سازی تیخونوف را به عنوان راه‌حلی برای رویارویی با مسائل بدشرطیده معرفی نمود. البته او این روش را علاوه بر مسائل بدشرطیده برای موضوعات بسیاری نظیر معادلات انتگرال، توپولوژی، تحلیل عملکردی و فیزیک ریاضی به کار برد. پس از او دیوید فیلیس<sup>۲۳</sup> به طور گسترده‌ای از این روش استفاده کرد. بدین سان در برخی از منابع، از این روش، تحت عنوان منظم‌سازی تیخونوف-فیلیس یاد شده است. نخستین بار هورل [۱۶] رگرسیون ستیغی را به عنوان یک روش انتخاب متغیر در آمار، به رسمیت شناخت.

<sup>17</sup>Penalty function

<sup>18</sup>Regularization parameter

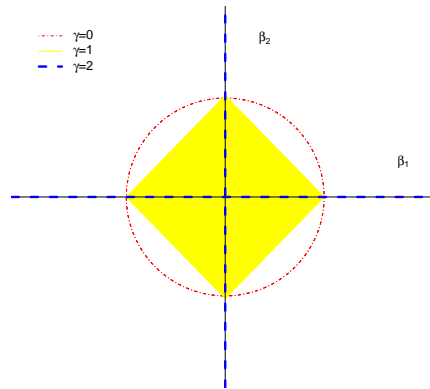
<sup>19</sup>Penalty parameter

<sup>20</sup>Square error loss function

<sup>21</sup>Bridge regression

<sup>22</sup>Andrey Nikolayevich Tikhonov

<sup>23</sup>David L. Phillips



شکل ۱: ناحیه محدودیت بریج.

به همین جهت عناصر قطری بزرگ‌تر یا مساوی صفر هستند. اکنون با نمایش  $\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p$  به صورت  $\mathbf{V} \mathbf{D} \mathbf{V}^T + \mathbf{V} \lambda \mathbf{I}_p \mathbf{V}^T$  واضح است که  $\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p$  قطری متعامد است. بنابراین

$$\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p = \mathbf{V} (\mathbf{D} + \lambda \mathbf{I}_p) \mathbf{V}^T$$

به طوری که تمام مقادیر ویژه آن (با توجه به این که  $\lambda > 0$  است) بزرگ‌تر از صفرند و زمانی که تمام مقادیر ویژه مخالف صفر باشند، ماتریس مربوطه وارون‌پذیر است. مسئله بهینه‌سازی روش ستیغی به صورت زیر است:

$$\begin{aligned} & \min_{\beta \in \mathbb{R}^p} \left\{ (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) + \lambda \beta^T \beta \right\} \\ & = \min_{\beta \in \mathbb{R}^p} \left\{ \sum_{i=1}^n (y_i - \sum_{j=1}^p x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^p \beta_j^2 \right\} \end{aligned}$$

که در آن تابع تاوان روش ستیغی  $\sum_{j=1}^p \beta_j^2$  می‌باشد. به  $\lambda$  پارامتر ستیغی<sup>۲۶</sup> نیز گفته می‌شود و برخی آن را با  $k$  نیز نمایش می‌دهند. مسئله بهینه‌سازی روش ستیغی محذب است و بدین سان مورد علاقه محققین است.

برای یافتن درک بصری از روش ستیغی با فرض وجود تنها دو متغیر توضیحی در مدل رگرسیونی ناحیه محدودیت به صورت

$$\beta_1^2 + \beta_2^2 \leq t \quad (5)$$

است. عبارت (۵) نشان‌دهنده معادله دایره‌ای به مرکز مبدأ مختصات و شعاع  $\sqrt{t}$  می‌باشد. بنابراین ناحیه محدودیت یک دایره است.

می‌شود که یکی از معایب بزرگ این روش است. بزرگ‌ترین مزیت روش ستیغی وارون‌پذیری همیشگی ماتریس  $\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p$  است. سؤال مهمی که در این قسمت پیش می‌آید، این است که چرا در روش ستیغی، وارون همیشه قابل محاسبه است؟ برای پاسخ به این سؤال از تجزیه مقادیر منفرد<sup>۲۴</sup> ماتریس‌ها استفاده می‌کنیم. هر ماتریسی حتی ماتریس‌های غیرمربعی دارای تجزیه مقادیر منفرد هستند. این روش یک ماتریس را به سه ماتریس دیگر تجزیه می‌کند. به عنوان مثال ماتریس  $\mathbf{G}$  را می‌توان به صورت  $\mathbf{G} = \mathbf{U} \mathbf{D} \mathbf{V}^T$  تجزیه نمود. به طوری که  $\mathbf{U}$  و  $\mathbf{V}$  ماتریس‌های متعامد و  $\mathbf{D}$  ماتریس قطری است. البته اگر ماتریس  $\mathbf{G}$  متقارن باشد، آنگاه  $\mathbf{V} = \mathbf{U}^T$  به صورت قطری متعامد<sup>۲۵</sup> است [۲۶]. به همین ترتیب ماتریس  $\mathbf{X}^T \mathbf{X}$  نیز دارای تجزیه SVD است. از طرفی این ماتریس متقارن است در نتیجه به صورت متعامد قطری‌پذیر است. بنابراین می‌توان این ماتریس را به صورت

$$\mathbf{X}^T \mathbf{X} = \mathbf{V} \mathbf{D} \mathbf{V}^T$$

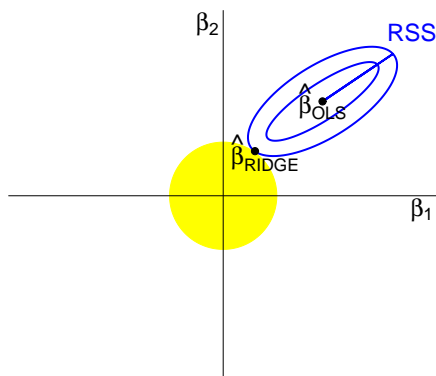
نوشت، که در آن  $\mathbf{D}$  یک ماتریس قطری است که عناصر روی قطر آن مقادیر ویژه‌ی ماتریس  $\mathbf{X}^T \mathbf{X}$  بوده و  $\mathbf{V}$  ماتریس متعامد است که ستون‌های آن بردارهای ویژه  $\mathbf{X}^T \mathbf{X}$  است. از سوی دیگر ماتریس  $\mathbf{X}^T \mathbf{X}$  نیمه معین مثبت است زیرا

$$\forall \mathbf{z} \neq \mathbf{0}, \quad \mathbf{z}^T \mathbf{X}^T \mathbf{X} \mathbf{z} = (\mathbf{Xz})^T (\mathbf{Xz}) = \|\mathbf{Xz}\|^2 \geq 0$$

<sup>24</sup>Singular value decomposition(SVD)

<sup>25</sup>Orthogonal diagonalization

<sup>26</sup>Ridge parameter



شکل ۲: ناحیه محدودیت ستیغی.

نامنفی برایمن که گروتی  $3^\circ$  نام دارد، به صورت

$$\sum_{i=1}^n \left( y_i - \sum_{j=1}^p c_j \hat{\beta}_j x_{ji} \right)^2 \quad c_j \geq 0, \quad \sum_j c_j \leq s$$

است که در آن  $\hat{\beta}_j$  برآوردگر اولیه کمترین توان‌های دوم است. با کاهش  $s$  (پارامتر منظم‌سازی) گروتی محدود می‌شود. روش گروتی برخی از متغیرها را حذف می‌کند و مابقی آن‌ها را منقبض می‌کند. بنابراین این روش مدل‌های قابل تفسیر ارائه می‌دهد. این روش نسبتاً پایدار است و همچنین مقیاس پایاست [۲۳]. مطالعات برایمن [۹] روی داده‌های واقعی و داده‌های شبیه‌سازی شده نشان داد که خطای پیش‌گویی گروتی در مقایسه با رگرسیون بهترین زیرمجموعه کمتر است و در رقابتی پایایی با رگرسیون ستیغی است، به جز زمانی که مدل واقعی ضرایب بسیار کوچک غیر صفر دارد. وابستگی گروتی به علامت و مقدار برآورد کمترین توان‌های دوم نقص بزرگی است که بی‌اعتمادی محققان به این روش را به ارمغان آورده است. بدین جهت می‌توان گفت زمانی که با داده‌های بعد بالا مواجه هستیم و یا زمانی که متغیرها بسیار همبسته‌اند، برآوردهای کمترین توان‌های دوم عملکرد ضعیفی دارد و به تبع آن گروتی نیز عملکرد ضعیفی داشته و مورد اطمینان نیست. این در حالی است که لاسو برخلاف گروتی از استفاده مستقیم از برآوردهای کمترین توان‌های دوم اجتناب می‌کند. مسئله بهینه‌سازی لاسو به صورت

$$\min_{\beta \in \mathbb{R}^p} \left\{ \sum_{i=1}^n (y_i - \sum_{j=1}^p x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^p |\beta_j| \right\},$$

در شکل ۲ دایره زرد رنگ ناحیه توان‌ها است. محل برخورد RSS (مجموع توان‌های دوم باقی‌مانده‌ها) روش کمترین توان‌های دوم با دایره، مختصات برآوردگر ستیغی را نمایش می‌دهد. با عمود کردن این نقطه بر محور افقی،  $\hat{\beta}_1$  و با عمود کردن آن بر محور عمودی،  $\hat{\beta}_2$  حاصل می‌شود. روش ستیغی فرآیندی پیوسته است که ضرایب تولیدی آن پایدار هستند اما برآوردهایی اریب ارائه می‌دهد که پذیرفتن مقداری اریبی در مقابل یافتن پاسخی قابل اطمینان‌تر معقول به نظر می‌رسد.

روش ستیغی ضرایب را به سمت صفر منقبض می‌کند، هرچه مقدار پارامتر ستیغی بیشتر باشد ضرایب برآورد شده انقباض بیشتری داشته و کوچک‌تر می‌شوند ولی هیچ‌گاه حتی در بی‌نهایت نیز به صفر نمی‌رسند این امر به سبب خاصیت دایره (ناحیه توان‌ها) است (ناحیه توان‌ها) رخ می‌دهد. البته به ازای  $p > 2$  نیز هیچ‌گاه ضرایب برآورد شده به صفر نمی‌رسند. بدین‌سان به کارگیری این روش برای داده‌های بعد بالا، مدل‌های قابل تفسیر ارائه نمی‌دهد و این روش تنگ نیست.

## ۶.۲ رگرسیون لاسو

رگرسیون لاسو  $27$  توسط رابرت تیبشیرانی  $28$  در سال ۱۹۹۶ پیشنهاد شد. ایده اولیه لاسو برگرفته از پیشنهاد جالب لئو برایمن  $29$  بود [۲۳]. مسئله بهینه‌سازی

<sup>27</sup>Lasso

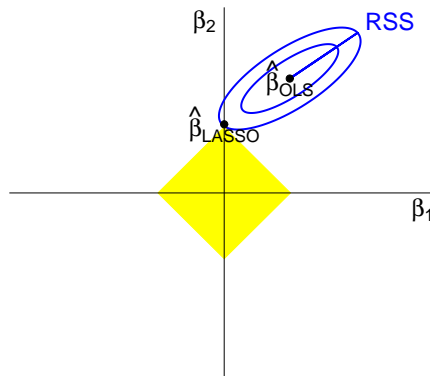
<sup>28</sup>Robert Tibshirani

<sup>29</sup>Leo Breiman

<sup>30</sup>Garrote

در مدل را کاهش می‌دهد. ناپایداری در مدل، یعنی با ایجاد کوچک‌ترین تغییراتی نتیجه انتخاب مدل بسیار متفاوت خواهد شد و دقت نتایج کاهش می‌یابد. لاسو به معنی عملگر انتخاب و کمترین قدرمطلق انقباضی<sup>۳۱</sup> است. لاسو برخی ضرایب را منقبض و بقیه را صفر می‌کند. بدین سبب هر دو ویژگی خوب رگرسیون انتخاب زیرمجموعه و ستیگی را حفظ می‌کند.

است. تیشیرانی [۲۳] برای حل مسئله بهینه‌سازی بالا یک مسئله برنامه‌ریزی درجه دوم با محدودیت نابربری خطی به کار برده است. این محدودیت به طور طبیعی تمایل به تولید ضرایبی که صفر هستند، دارد و از این رو مدل‌های قابل تفسیر ارائه می‌دهد. روش لاسو برآوردگرهای پایدار، اریب و پیوسته تولید می‌کند [۲۳]. یعنی برآوردگر تولیدشده روی داده‌ها پیوسته است و ناپایداری



شکل ۳: ناحیه محدودیت لاسو.

## ۷.۲ روش شبکه ارتجاعی

برای یافتن درک بصری از روش لاسو با فرض وجود تنها دو متغیر توضیحی در مدل رگرسیونی، ناحیه محدودیت به صورت

$$|\beta_1| + |\beta_2| \leq t$$

این روش که ترکیبی از روش لاسو و ستیگی است توسط ژو و هستی [۲۷] معرفی شد. آن‌ها این روش را برای تعدیل هم‌زمان مشکلات هم‌خطی چندگانه و تفسیر ناپذیر بودن مدل، مطرح کردند.

است که با رسم آن، ناحیه محدودیت یک لوزی است. در شکل ۳ لوزی زرد رنگ ناحیه تاوان است. محل برخورد RSS روش کمترین توان‌های دوم با لوزی، مختصات برآوردگر لاسو را نمایش می‌دهد. با عمود کردن این نقطه بر محور افقی،  $\beta_1$  و با عمود کردن آن بر محور عمودی،  $\beta_2$  حاصل می‌شود. اگر محل برخورد دقیقاً روی یکی از گوشه‌های لوزی باشد، (همانند شکل ۳) تنها یک متغیر انتخاب شده و دیگری صفر می‌شود. همان‌طور که ذکر شد، در روش لاسو به سبب نوع ناحیه تاوان (لوزی و گوشه داشتن لوزی) برخی از متغیرها دقیقاً برابر صفر می‌شوند. از منظر دیگر می‌توان استدلال کرد که چون تابع تاوان روش لاسو در نقطه صفر مشتق ناپذیر است، مقادیر بزرگ‌تری از  $\lambda$  منجر به حذف متغیرهای پیشگوی کم تأثیر در مدل می‌گردد و در نتیجه عمل برآورد پارامتر و انتخاب متغیر به صورت هم‌زمان انجام می‌گیرد [۳].

مسئله بهینه‌سازی به صورت

$$\min_{\beta \in \mathbb{R}^p} \left\{ \sum_{i=1}^n (y_i - \sum_{j=1}^p x_{ij} \beta_j)^2 + \lambda_1 \sum_{j=1}^p |\beta_j| + \lambda_2 \sum_{j=1}^p \beta_j^2 \right\}$$

روش شبکه ارتجاعی<sup>۳۲</sup> دو پارامتر تاوان  $\lambda_1$  و  $\lambda_2$  را داراست که طبیعتاً یافتن مقدار مناسب برای دو پارامتر تاوان این روش نسبت به روش‌هایی که تنها یک پارامتر تاوان دارند، امری به مراتب دشوارتر و سخت‌تر خواهد بود و این یکی از معایب این روش محسوب می‌شود. تابع تاوان این روش، محدب است.

<sup>۳۱</sup>Least Absolute Shrinkage and Selection Operator (LASSO)

<sup>۳۲</sup>Elastic Net



## ۸.۲ روش اسکد

بودن مسئله بهینه‌سازی آن است. فن و لی [۱۲] مقدار  $a = 3/7$  را به‌عنوان مقداری مناسب برای حل مسئله بهینه‌سازی این روش پیشنهاد نمودند. در شکل ۴ تابع تاوان روش لاسو، ستیخی و اسکد به‌طور هم‌زمان رسم شده‌اند.

فن و لی [۱۲] استفاده از روش اسکد<sup>۳۳</sup> را که تابع تاوان آن توسط [۱۳] ارائه شده بود، پیشنهاد کردند. تابع تاوان پیشنهادی آن‌ها به‌صورت زیر است:

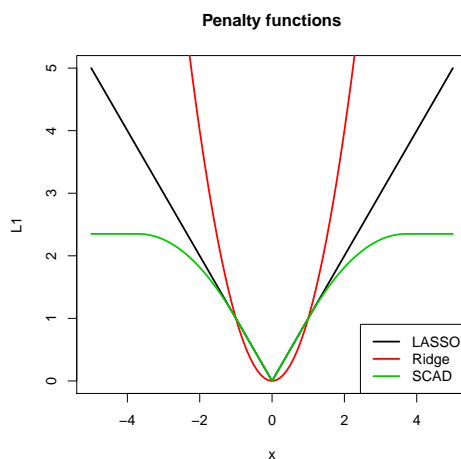
$$P_{\lambda}(\beta_j) = \begin{cases} \lambda |\beta_j| & \text{اگر } |\beta_j| \leq \lambda; \\ -\left(\frac{|\beta_j|^2 - 2a\lambda|\beta_j| + \lambda^2}{2(a-1)}\right) & \text{اگر } \lambda < |\beta_j| \leq a\lambda; \\ \frac{(a+1)\lambda^2}{4} & \text{اگر } |\beta_j| > a\lambda \end{cases}$$

## ۹.۲ روش اسکار

باندل و ریچ [۸] روش اسکار<sup>۳۴</sup> را منتشر کردند. این روش در صورت وجود هم‌خطی عملکرد خوبی دارد. ناحیه تاوان این روش به‌صورت زیر است.

$$\min_{\beta \in \mathbb{R}^p} \left\{ \sum_{i=1}^n (y_i - \sum_{j=1}^p \mathbf{X}_{ij} \beta_j)^2 + \lambda \left( \sum_{j=1}^p |\beta_j| + c \sum_{j < k} \max\{|\beta_j|, |\beta_k|\} \right) \right\} \quad c \geq 0.$$

$\lambda$  و  $a$  دو پارامتر تاوان این روش هستند که همین امر سبب دشواری این روش است. پارامتر  $a$  را با  $\gamma$  نیز نمایش می‌دهند. برآوردهای این روش ناریب و پیوسته هستند. در این روش نیز متغیرهای کم تأثیر از مدل خارج می‌شوند. به بیان ساده‌تر این روش تنک است. یکی از عیوب بزرگ این روش غیرمحدب



شکل ۴: تابع تاوان سه روش لاسو، ستیخی و اسکد.

خانواده‌ی برآوردهای کمترین توان‌های دوم تاوانیده<sup>۳۵</sup> خانواده‌ی وسیعی از برآوردها است که توصیف و ذکر تمامی این موارد در این مقاله نمی‌گنجد.

که در آن پارامتر  $c$  یک ثابت تنظیم‌کننده است. تاوان  $\sum_{j=1}^p |\beta_j|$  باعث تنگی و همچنین تاوان  $\sum_{j < k} \max\{|\beta_j|, |\beta_k|\}$  باعث تساوی ضرایب می‌شود. تابع هدف این روش نیز محدب است. هشت ضلعی موجود در شکل ۵ ناحیه تاوان این روش است.

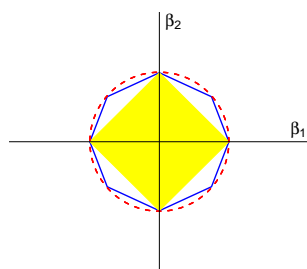
متغیرهای پراهمیت از مدل حذف شود به‌طوری که در بی‌نهایت تمامی متغیرها را از مدل خارج می‌کند. کاهش این پارامتر، نیز کاهش اریبی و افزایش واریانس را به ارمغان می‌آورد. همچنین باعث کوچکی ناحیه محدودیت می‌شود. بدین جهت ممکن است متغیرهای کم‌اهمیت زیادی در مدل باقی

یکی از معایب بزرگ روش‌های تاوانیده، وابستگی آن‌ها به پارامتر تاوان است. افزایش پارامتر تاوان، افزایش اریبی و کاهش واریانس را به همراه دارد. همچنین باعث بزرگی ناحیه محدودیت شده که این امر در روش‌های تنک منجر به حذف متغیرهای فراوانی می‌شود. بدین‌سان حتی ممکن است

<sup>33</sup>Smoothly Clipped Absolute Deviation (SCAD)

<sup>34</sup>Octagonal shrinkage and clustering algorithm for regression (OSCAR)

<sup>35</sup>The penalized least squares (PLS) family



شکل ۵: ناحیه محدودیت اسکار.

یعنی

$$E(\hat{\beta}) = \beta$$

۳. پیوستگی: تابع توان موردنظر منجر به ایجاد روشی پایدار و برآوردگرهایی پیوسته شود.

یکی دیگر از موارد بسیار مهم در روش توانانیده که رویه انتخاب متغیر دارند، این است که بتوانند زیرمجموعه درست و واقعی متغیرها را شناسایی کنند. این خاصیت الهام‌بخش، تحت عنوان خاصیت پیشگویی<sup>۳۹</sup> شناخته می‌شود که برای نخستین بار توسط فن و لی [۱۲] مطرح شد.

لازم به ذکر است که بررسی این خاصیت در عمل برای داده‌های شبیه‌سازی شده امکان‌پذیر است.

در میان تمامی روش‌های مطرح‌شده، روش‌هایی که تابع توان آن‌ها محذب است، موردپسندتر است. چراکه اثبات می‌شود در مسائل بهینه‌سازی محذب<sup>۴۰</sup>، مینیمم (ماکزیمم) نسبی با مینیمم (ماکزیمم) مطلق معادل است.

## ۱۱.۲ روش توانانیده اصلاح‌شده

همان‌طور که مطرح شد روش‌های توانانیده دارای معایبی نظیر انتخاب مقدار بهینه‌ی پارامتر توان و ... است که نقش مهمی در نتیجه‌ی برآورد دارد، به این معنی که مقادیر بزرگ پارامتر توان منجر به بی‌ثباتی قابل توجهی می‌شود. انتخاب مقدار پارامتر توان باوجود تمامی تحقیقات انجام‌شده و روش‌های متنوع پیشنهادشده توسط محققین رشته‌های آمار و ریاضی هنوز به‌طور کامل حل‌نشده است و همچنان عرصه برای تلاش بیشتر و معرفی روش‌های دقیق‌تر

بمانند. با انتخاب صفر برای این پارامتر روش توانانیده دقیقاً معادل روش کمترین توان‌های دوم خواهد بود. بدین منظور یکی از نکات بسیار مهم در روش توانانیده تعیین مقدار مناسبی برای پارامتر توان است. منظور از مقدار مناسب، مقداری است که توازنی بین اریبی و واریانس برقرار نماید، به‌طوری‌که RSS مدل را به کمترین مقدار ممکن برساند. برای تعیین این مقدار بهینه روش‌هایی نظیر معیار اطلاع آکائیکه<sup>۳۶</sup>، اطلاع بیز<sup>۳۷</sup>، آکائیکه تصحیح‌شده،  $C_p$  مالوس، اعتبارسنجی<sup>۳۸</sup> و ... وجود دارد. اخیراً رایج‌ترین روش در میان محققان برای محاسبه‌ی این پارامتر، روش اعتبارسنجی است. دلیل محبوبیت این روش، محاسبه‌ی هم‌زمان مقدار بهینه و آزمون میزان دقت آن می‌باشد.

## ۱۰.۲ تابع توان

در بخش قبل برخی از روش‌های توانانیده به‌اختصار بیان شد. اما واقعیت این است که توابع توان فراوانی وجود دارد که هر کدام منجر به ایجاد روش جدیدی می‌شوند. سؤال مهمی که در این قسمت پیش می‌آید، این است که کدام از این توابع توان موردپسند هستند و یا این که ویژگی‌های مطلوب یک تابع توان چیست؟ فن و لی [۱۲] با انتشار مقاله‌ای به این سؤال‌ها پاسخ دادند. از نظر آن‌ها یک تابع توان مطلوب باید هم‌زمان سه ویژگی مهم که به شرح ذیل است، را دارا باشد.

۱. تنکی: تابع توان موردنظر باید ضرایبی که مقدارشان کوچک است را برابر صفر قرار دهد. همین امر منجر به کاهش پیچیدگی و افزایش تفسیرپذیری مدل خواهد شد.
۲. نارایی: تابع توان موردنظر باید منجر به ایجاد برآوردی نارایب شود،

<sup>36</sup>Akaike information criterion (AIC)

<sup>37</sup>Bayesian information criterion (BIC)

<sup>38</sup>Cross-validation(CV)

<sup>39</sup>Oracle property

<sup>40</sup>Convex optimization

و برای تولید خطا از

$$\varepsilon \sim N(0, 1)$$

استفاده شد. در برخی از روش‌های به کاررفته در این مثال نیازمند استفاده از روش اعتبارسنجی متقابل هستیم. در تمام این روش‌ها مجموعه آموزش ۷۰٪ و مجموعه آزمون ۳۰٪ در نظر گرفته شده است. در گام اول از روش مؤلفه‌های اصلی استفاده می‌شود که نتایج آن در جدول ۱ خلاصه شده است. با توجه به جدول ۱ به ازای ۱۷ مؤلفه ۶۲٪ تغییرات، به ازای ۲۱ مؤلفه ۷۱٪ تغییرات، به ازای ۲۶ مؤلفه ۸۰٪ تغییرات، به ازای ۳۴ مؤلفه ۹۰٪ تغییرات، قابل توجیه بوده و از مؤلفه چهل و هشتم به بعد ۱۰۰ درصد تغییرات توسط واریانس کل توجیه می‌شود. نمودار بازو و لگاریتم آن در شکل ۶ نمایش داده شده است. تعداد مؤلفه‌ها با این نمودار و نمودار لگاریتم ۴۹ مؤلفه تعیین می‌شود. در ادامه از روش‌های تاوانیده استفاده می‌کنیم در شکل ۷ و ۸ نمودارهای مربوط به اعتبارسنجی و برآورد ضرایب به ازای پارامترهای تاوانیده مختلف با روش ستیگی و لاسو مشهود است. در نمودار اعتبارسنجی روش لاسو به ازای مقدار ۰/۸۰۹۷ مقدار MSE را کمینه کرده و در روش لاسو مقدار آن به ازای ۲۳/۱۰۱۳ کمینه شده است. در قسمت فوقانی نمودارهای سمت راست شکل‌های ۷ و ۸ تعداد ضرایب غیر صفر موجود به ازای مقادیر مختلف  $\lambda$  نمایش داده شده است که واضح است که با افزایش میزان  $\lambda$  تعداد ضرایب غیر صفر روش لاسو کاهش می‌یابد درحالی که در روش ستیگی برآورد ضرایب بسیار به صفر نزدیک می‌شوند ولی دقیقاً معادل با صفر نخواهند بود. نمودار اعتبارسنجی برای روش شبکه ارتجاعی در شکل ۹ نمایش داده شده است. نمودارهای برآورد ضرایب به روش شبکه ارتجاعی به ازای مقادیر مختلف پارامتر  $\alpha$  و بردار پارامتر دلخواه یکسان در شکل ۱۲ رسم شده است. با توجه به این شکل‌ها بدیهی است که با افزایش میزان  $\alpha$  متغیرهای بیشتری از مدل خارج می‌شوند. نمودارهای برآورد ضرایب به روش SCAD به ازای مقادیر مختلف پارامتر  $a$  و بردار پارامتر دلخواه یکسان در شکل ۱۰ رسم شده است. نمودار اعتبارسنجی متقابل برای یافت مقدار بهینه پارامتر  $\lambda$  به ازای مقدار پیشنهادی فن و لی [۱۲] برای پارامتر  $a$  در شکل ۱۱ نشان داده شده است. در جدول ۲ نتایج روش‌های مورد استفاده گزارش شده است.

می‌شود، یکی از ویتامین‌های B است که همه محلول در آب هستند. از مهم‌ترین ویتامین‌های لازم برای حیات موجود زنده است. دلیل انتخاب این نام، رنگ زرد این ویتامین است که ناشی از وجود حلقه فلاوینی موجود در ساختمان آن می‌باشد. ریوفلاوین به‌طور طبیعی در برخی غذاها وجود دارد، به

به‌منظور انتخاب بهینه‌ترین مقدار وجود دارد. زیرا برآورد به‌دست‌آمده در روش‌های تاوانیده وابستگی زیادی به پارامتر تاوان مجهول دارد.

در واقع مشکل بنیادینی که داده‌ی بعد بالا برای ما به ارمغان می‌آورد عدم وارون‌پذیری ماتریس  $\mathbf{X}^T \mathbf{X}$  به سبب بدشرطیده<sup>۴۱</sup> بودن آن است. ماتریس بدشرطیده دارای عدد شرطی<sup>۴۲</sup> بزرگی است. به همین شکل برای مقابله با این مشکل تنها کافی است ماتریس  $\mathbf{X}^T \mathbf{X}$  با کاهش عدد شرطی به یک ماتریس خوش شرطیده<sup>۴۳</sup> تبدیل شود. به همین سبب کافی است که ضریبی مثبت از عدد شرطی ماتریس  $\mathbf{X}^T \mathbf{X}$  به تابع هدف به‌عنوان عبارت تاوان، به‌منظور کنترل عدد حالت ماتریس افزوده می‌شود. بنابراین با استفاده از موارد ذکر شده و استفاده از روش پیشنهادی روزبه و همکاران [۲۲] می‌توان نوع اصلاح‌شده مسئله بهینه‌سازی (۴) را به‌صورت

$$\min_{\beta \in \mathbb{R}^p} \left\{ (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) + M \mathcal{H}(\mathbf{X}^T \mathbf{X}) \right\},$$

نوشت، که در آن  $M \geq 0$  پارامتر تاوان نامیده شده و  $\mathcal{H}(\mathbf{X}^T \mathbf{X})$  عدد شرطی ماتریس  $\mathbf{X}^T \mathbf{X}$  است. اعمال این روش موجب رهایی از دردسرهای روش‌های تاوانیده می‌شود، اما متأسفانه تنگ نیست. این روش برای مقابله با هم‌خطی نیز سودمند خواهد بود.

### ۳ مطالعه شبیه‌سازی

در این بخش در یک مطالعه شبیه‌سازی به بررسی و کاربرد روش‌های بیان شده می‌پردازیم. در این مطالعه نمونه‌ای با اندازه ۵۰ به همراه ۱۰۰ متغیر توضیحی از مدل

$$\mathbf{y} = \mathbf{X}\beta + \varepsilon$$

تولید شده‌اند. برای تولید بردار پارامترها از

$$\beta_1 = (-1, 5, 2, 2, 5, 4, -3, 5)^T, \quad \beta_2 \sim N(0, 0, 1),$$

برای تولید ماتریس طرح از

$$\mathbf{X} \sim N_p(\mu, \mathbf{I}), \quad \mu = (0, \dots, 0)^T_{p \times 1}$$

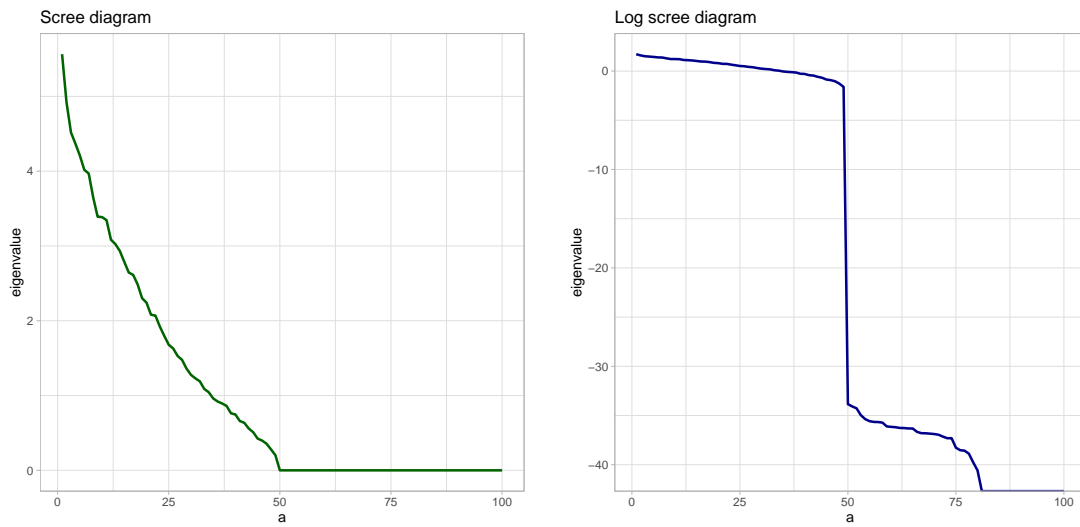
### ۴ مطالعه‌ی داده واقعی

در این بخش به‌منظور بررسی یک مجموعه داده واقعی با بعد بالا از داده‌های ریوفلاوین استفاده می‌شود. ریوفلاوین که به‌عنوان ویتامین B2 نیز شناخته

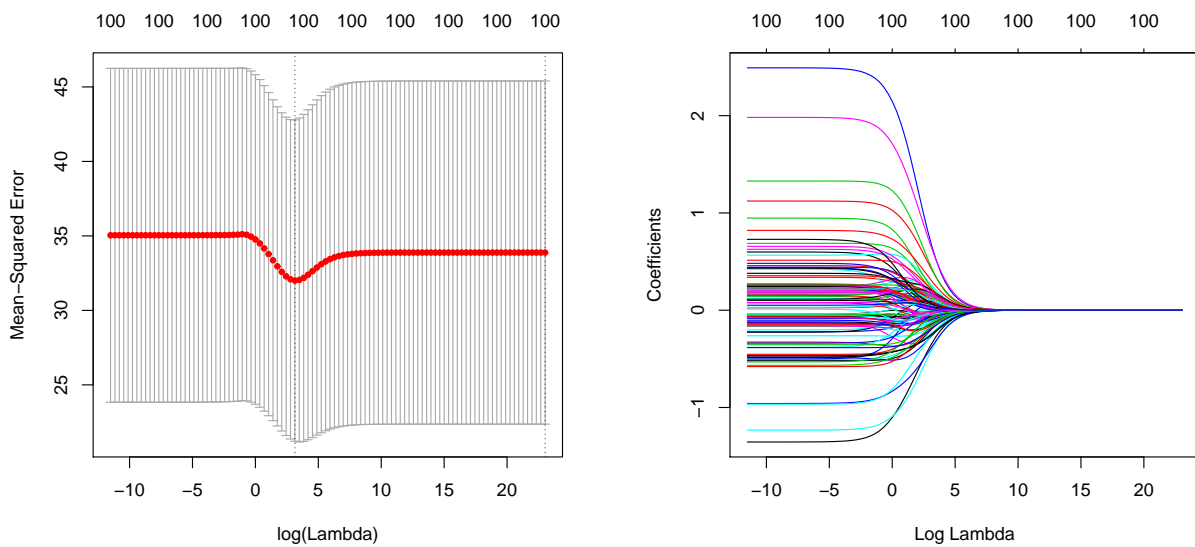
<sup>41</sup>Ill-condition

<sup>42</sup>Condition number

<sup>43</sup>Well-condition



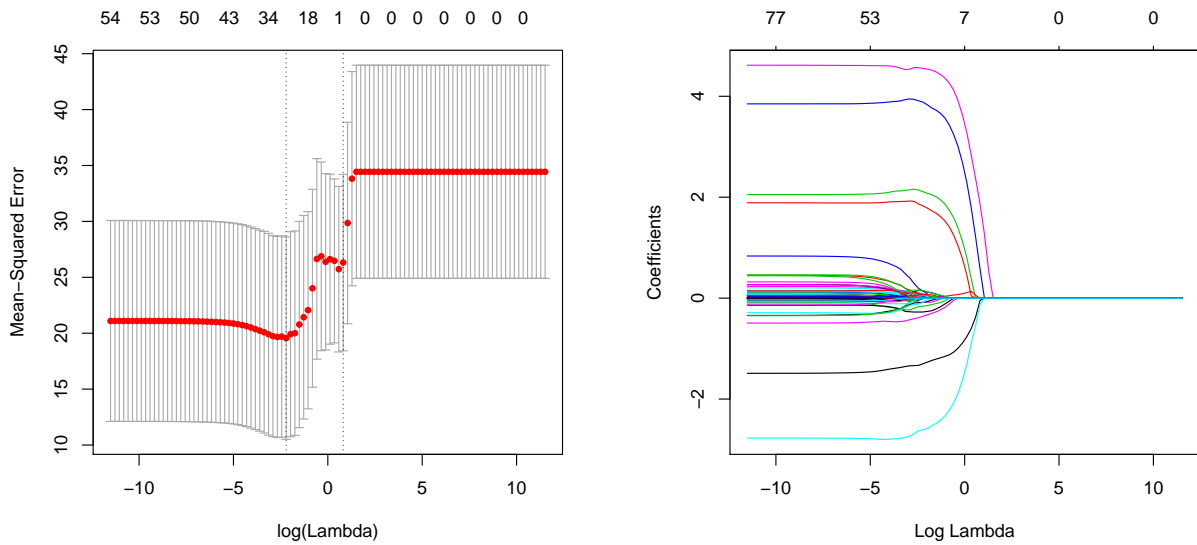
شکل ۶: نمودار بازو و لگاریتم آن برای داده‌های شبیه‌سازی شده.



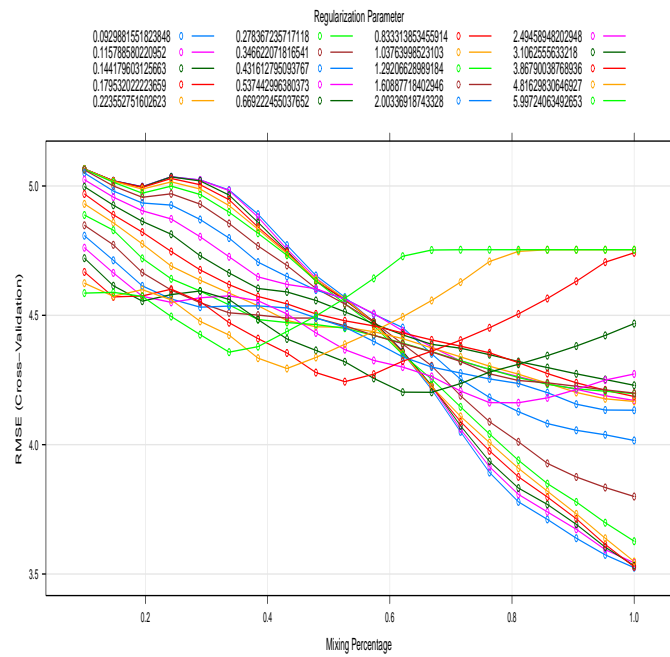
شکل ۷: نمودار اعتبارسنجی و برآورد ضرایب داده‌های شبیه‌سازی شده به ازای پارامترهای توان مختلف به روش ستیخی.

جدول ۱: نتایج روش مؤلفه‌های اصلی برای داده‌های شبیه‌سازی شده.

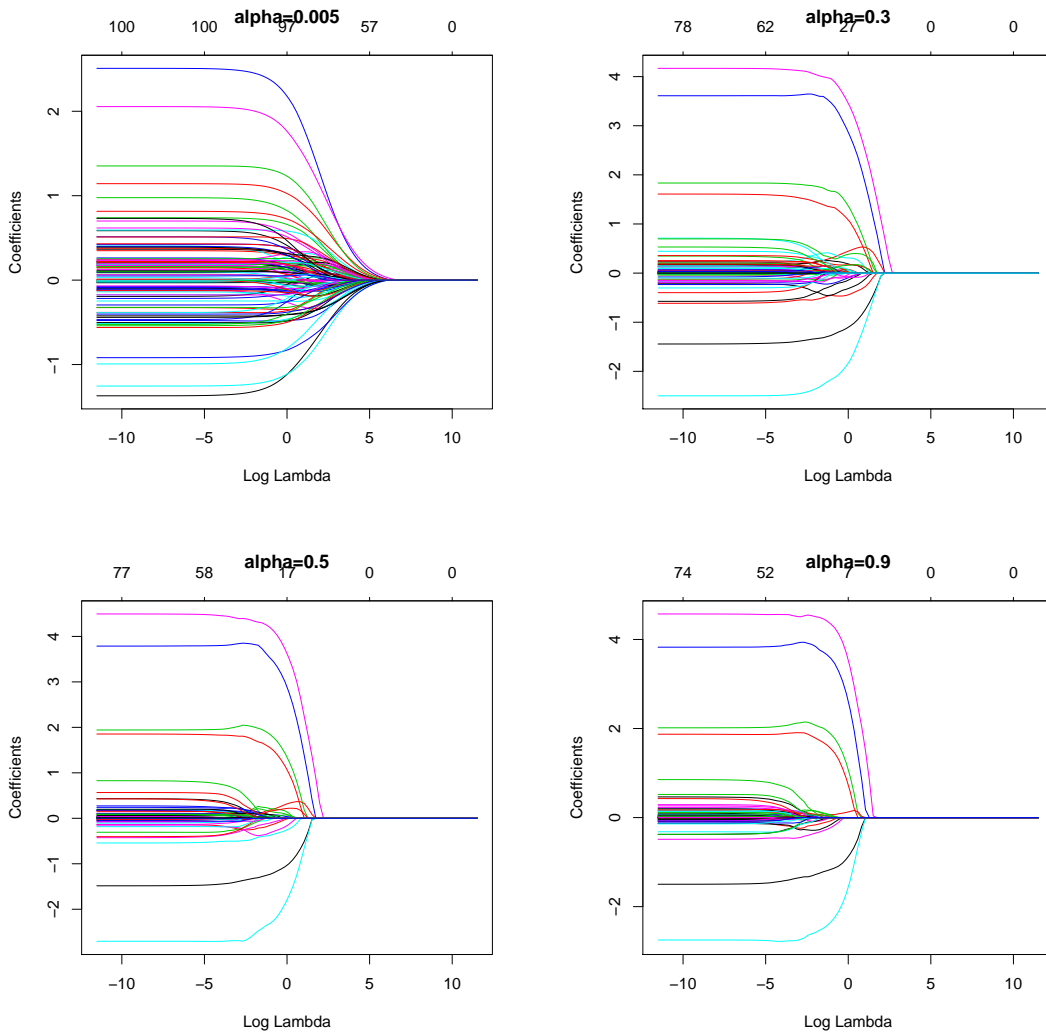
مؤلفه	نسبت واریانس	نسبت تجمعی واریانس	مؤلفه	نسبت واریانس	نسبت تجمعی واریانس
مؤلفه اول	۰/۵۵۶	۰/۵۵۶	مؤلفه پنجاه و یکم	۰/۵۵۶	۰/۵۵۶
مؤلفه دوم	۰/۴۹۱	۰/۱۰۴۸	مؤلفه پنجاه و دوم	۰/۱۰۴۸	۰/۱۰۴۸
مؤلفه سوم	۰/۴۵۱	۰/۱۵۰۰	مؤلفه پنجاه و سوم	۰/۱۵۰۰	۰/۱۵۰۰
مؤلفه چهارم	۰/۴۳۶	۰/۱۹۳۶	مؤلفه پنجاه و چهارم	۰/۱۹۳۶	۰/۱۹۳۶
مؤلفه پنجم	۰/۴۲۰	۰/۲۳۵۷	مؤلفه پنجاه و پنجم	۰/۲۳۵۷	۰/۲۳۵۷
مؤلفه ششم	۰/۴۰۱	۰/۲۷۵۹	مؤلفه پنجاه و ششم	۰/۲۷۵۹	۰/۲۷۵۹
مؤلفه هفتم	۰/۳۹۶	۰/۳۱۵۶	مؤلفه پنجاه و هفتم	۰/۳۱۵۶	۰/۳۱۵۶
مؤلفه هشتم	۰/۳۶۴	۰/۳۵۲۱	مؤلفه پنجاه و هشتم	۰/۳۵۲۱	۰/۳۵۲۱
مؤلفه نهم	۰/۳۳۹	۰/۳۸۶۰	مؤلفه پنجاه و نهم	۰/۳۸۶۰	۰/۳۸۶۰
مؤلفه دهم	۰/۳۳۸	۰/۴۱۹۹	مؤلفه شصتم	۰/۴۱۹۹	۰/۴۱۹۹
مؤلفه یازدهم	۰/۳۳۴	۰/۴۵۳۳	مؤلفه شصت و یکم	۰/۴۵۳۳	۰/۴۵۳۳
مؤلفه دوازدهم	۰/۳۰۸	۰/۴۸۴۲	مؤلفه شصت و دوم	۰/۴۸۴۲	۰/۴۸۴۲
مؤلفه سیزدهم	۰/۳۰۲	۰/۵۱۴۴	مؤلفه شصت و سوم	۰/۵۱۴۴	۰/۵۱۴۴
مؤلفه چهاردهم	۰/۲۹۳	۰/۵۴۳۷	مؤلفه شصت و چهارم	۰/۵۴۳۷	۰/۵۴۳۷
مؤلفه پانزدهم	۰/۲۷۸	۰/۵۷۱۶	مؤلفه شصت و پنجم	۰/۵۷۱۶	۰/۵۷۱۶
مؤلفه شانزدهم	۰/۲۶۴	۰/۵۹۸۱	مؤلفه شصت و ششم	۰/۵۹۸۱	۰/۵۹۸۱
مؤلفه هفدهم	۰/۲۶۱	۰/۶۲۴۲	مؤلفه شصت و هفتم	۰/۶۲۴۲	۰/۶۲۴۲
مؤلفه هجدهم	۰/۲۴۸	۰/۶۴۹۱	مؤلفه شصت و هشتم	۰/۶۴۹۱	۰/۶۴۹۱
مؤلفه نوزدهم	۰/۲۳۰	۰/۶۷۲۱	مؤلفه شصت و نهم	۰/۶۷۲۱	۰/۶۷۲۱
مؤلفه بیستم	۰/۲۲۴	۰/۶۹۴۵	مؤلفه هفتادم	۰/۶۹۴۵	۰/۶۹۴۵
مؤلفه بیست و یکم	۰/۲۰۸	۰/۷۱۵۳	مؤلفه هفتاد و یکم	۰/۷۱۵۳	۰/۷۱۵۳
مؤلفه بیست و دوم	۰/۲۰۶	۰/۷۳۶۰	مؤلفه هفتاد و دوم	۰/۷۳۶۰	۰/۷۳۶۰
مؤلفه بیست و سوم	۰/۱۹۲	۰/۷۵۵۲	مؤلفه هفتاد و سوم	۰/۷۵۵۲	۰/۷۵۵۲
مؤلفه بیست و چهارم	۰/۱۷۹	۰/۷۷۳۲	مؤلفه هفتاد و چهارم	۰/۷۷۳۲	۰/۷۷۳۲
مؤلفه بیست و پنجم	۰/۱۶۷	۰/۷۹۰۰	مؤلفه هفتاد و پنجم	۰/۷۹۰۰	۰/۷۹۰۰
مؤلفه بیست و ششم	۰/۱۶۲	۰/۸۰۶۳	مؤلفه هفتاد و ششم	۰/۸۰۶۳	۰/۸۰۶۳
مؤلفه بیست و هفتم	۰/۱۵۲	۰/۸۲۱۶	مؤلفه هفتاد و هفتم	۰/۸۲۱۶	۰/۸۲۱۶
مؤلفه بیست و هشتم	۰/۱۴۷	۰/۸۳۶۳	مؤلفه هفتاد و هشتم	۰/۸۳۶۳	۰/۸۳۶۳
مؤلفه بیست و نهم	۰/۱۳۵	۰/۸۴۹۹	مؤلفه هفتاد و نهم	۰/۸۴۹۹	۰/۸۴۹۹
مؤلفه سی ام	۰/۱۲۷	۰/۸۶۲۷	مؤلفه هشتادم	۰/۸۶۲۷	۰/۸۶۲۷
مؤلفه سی و یکم	۰/۱۲۲	۰/۸۷۵۰	مؤلفه هشتاد و یکم	۰/۸۷۵۰	۰/۸۷۵۰
مؤلفه سی و دوم	۰/۱۱۹	۰/۸۸۶۹	مؤلفه هشتاد و دوم	۰/۸۸۶۹	۰/۸۸۶۹
مؤلفه سی و سوم	۰/۱۰۸	۰/۸۹۷۸	مؤلفه هشتاد و سوم	۰/۸۹۷۸	۰/۸۹۷۸
مؤلفه سی و چهارم	۰/۱۰۴	۰/۹۰۸۲	مؤلفه هشتاد و چهارم	۰/۹۰۸۲	۰/۹۰۸۲
مؤلفه سی و پنجم	۰/۰۹۶	۰/۹۱۷۸	مؤلفه هشتاد و پنجم	۰/۹۱۷۸	۰/۹۱۷۸
مؤلفه سی و ششم	۰/۰۹۲	۰/۹۲۷۰	مؤلفه هشتاد و ششم	۰/۹۲۷۰	۰/۹۲۷۰
مؤلفه سی و هفتم	۰/۰۸۹	۰/۹۳۶۰	مؤلفه هشتاد و هفتم	۰/۹۳۶۰	۰/۹۳۶۰
مؤلفه سی و هشتم	۰/۰۸۶	۰/۹۴۴۶	مؤلفه هشتاد و هشتم	۰/۹۴۴۶	۰/۹۴۴۶
مؤلفه سی و نهم	۰/۰۷۶	۰/۹۵۲۲	مؤلفه هشتاد و نهم	۰/۹۵۲۲	۰/۹۵۲۲
مؤلفه چهل و یکم	۰/۰۷۴	۰/۹۵۹۷	مؤلفه نود	۰/۹۵۹۷	۰/۹۵۹۷
مؤلفه چهل و دوم	۰/۰۶۵	۰/۹۶۶۳	مؤلفه نود و یکم	۰/۹۶۶۳	۰/۹۶۶۳
مؤلفه چهل و سوم	۰/۰۶۳	۰/۹۷۲۶	مؤلفه نود و دوم	۰/۹۷۲۶	۰/۹۷۲۶
مؤلفه چهل و سوم	۰/۰۵۵	۰/۹۷۸۲	مؤلفه نود و سوم	۰/۹۷۸۲	۰/۹۷۸۲
مؤلفه چهل و چهارم	۰/۰۵۰	۰/۹۸۳۳	مؤلفه نود و چهارم	۰/۹۸۳۳	۰/۹۸۳۳
مؤلفه چهل و پنجم	۰/۰۴۲	۰/۹۸۷۶	مؤلفه نود و پنجم	۰/۹۸۷۶	۰/۹۸۷۶
مؤلفه چهل و ششم	۰/۰۳۹	۰/۹۹۱۶	مؤلفه نود و ششم	۰/۹۹۱۶	۰/۹۹۱۶
مؤلفه چهل و هفتم	۰/۰۳۵	۰/۹۹۵۱	مؤلفه نود و هفتم	۰/۹۹۵۱	۰/۹۹۵۱
مؤلفه چهل و هشتم	۰/۰۲۸	۰/۹۹۷۹	مؤلفه نود و هشتم	۰/۹۹۷۹	۰/۹۹۷۹
مؤلفه چهل و نهم	۰/۰۲۰	۱	مؤلفه نود و نهم	۱	۰/۰۲۰
مؤلفه پنجاهم	۲/۰۵۵ × ۱۰ <sup>-۱۷</sup>	۱	مؤلفه صدم	۱	۲/۰۵۵ × ۱۰ <sup>-۱۷</sup>



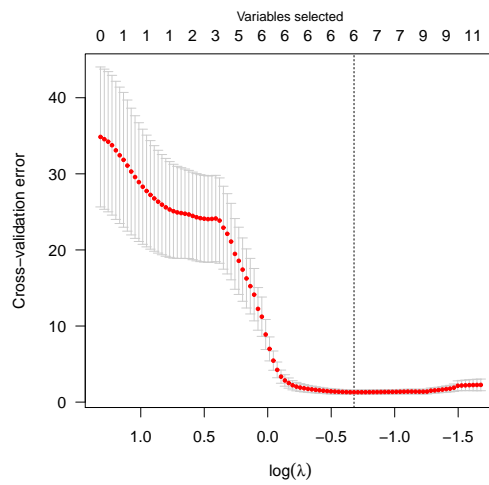
شکل ۸: نمودار اعتبارسنجی و برآورد ضرایب داده‌های شبیه‌سازی شده به ازای پارامترهای تاوان مختلف به روش لاسو.



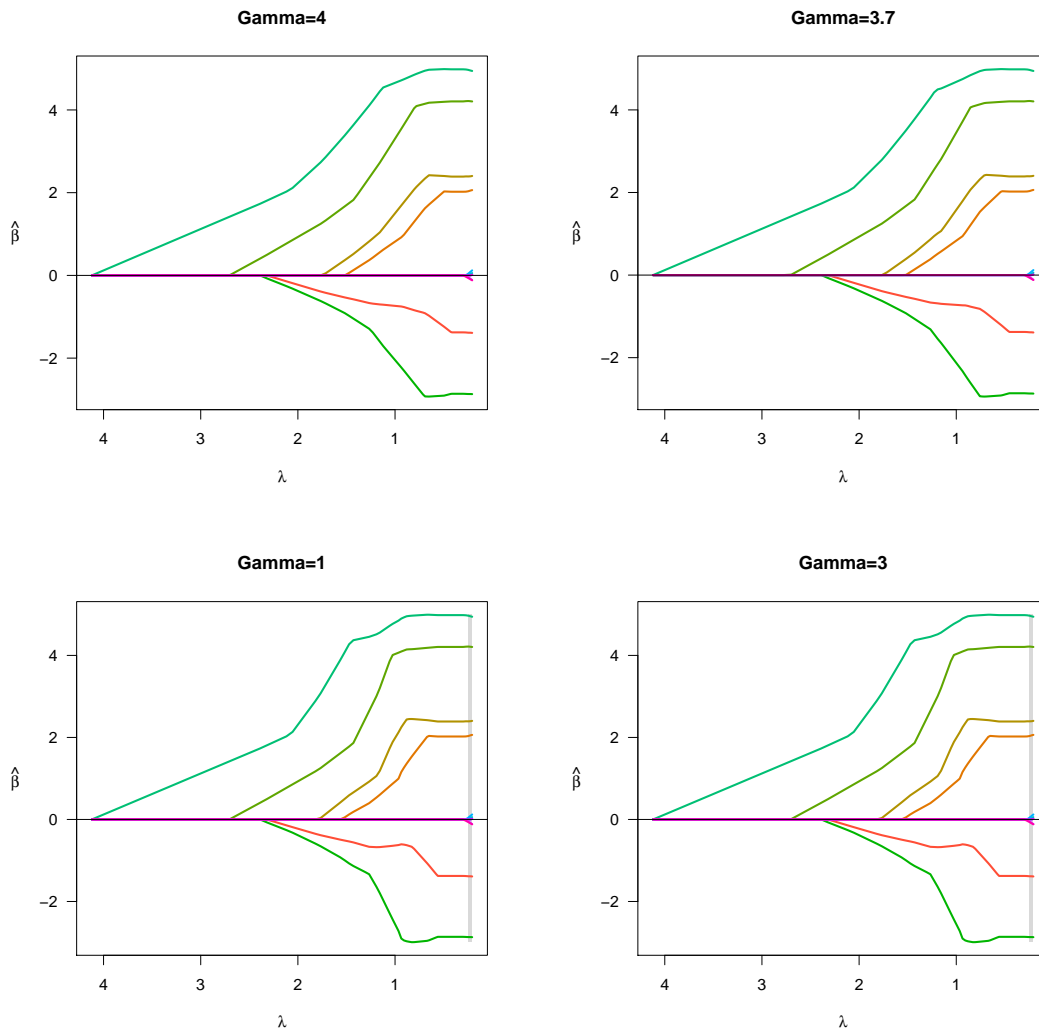
شکل ۹: نمودار اعتبارسنجی برای داده‌های شبیه‌سازی شده به روش شبکه ارتجاعی.



شکل ۱۰: نمودار برآورد ضرایب داده‌های شبیه‌سازی شده شبکه ارتجاعی.



شکل ۱۱: نمودار اعتبارسنجی داده‌های شبیه‌سازی شده برای پارامتر  $\lambda$  با روش SCAD.



شکل ۱۲: نمودار برآورد ضرایب داده‌های شبیه‌سازی شده به ازای پارامترهای توان مختلف به روش SCAD.

جدول ۲: جدول نتایج روش‌های توانیده برای داده‌های شبیه‌سازی شده

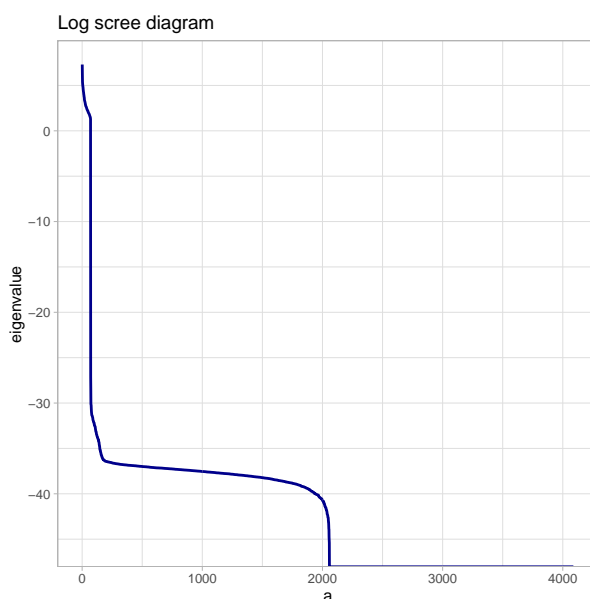
روش	مقدار بهینه پارامتر توان	تعداد ضرایب غیر صفر	مجموع مربعات خطا
ستیغی	۲۳/۱۰۱۳	۱۰۰	۱۹۳/۱۳۰۹
لاسو	۰/۱۰۹۷	۲۸	۲۲۳/۱۵۷
شبکه ارتجاعی	$\alpha = ۱, \lambda = ۰/۰۹۲۹۸$	۵۶	۵۳/۸۷۵۹
اسکد	$a = ۳/۷, \lambda = ۰/۵۰۶۷$	۶	۴۱/۹۱۴۰۵



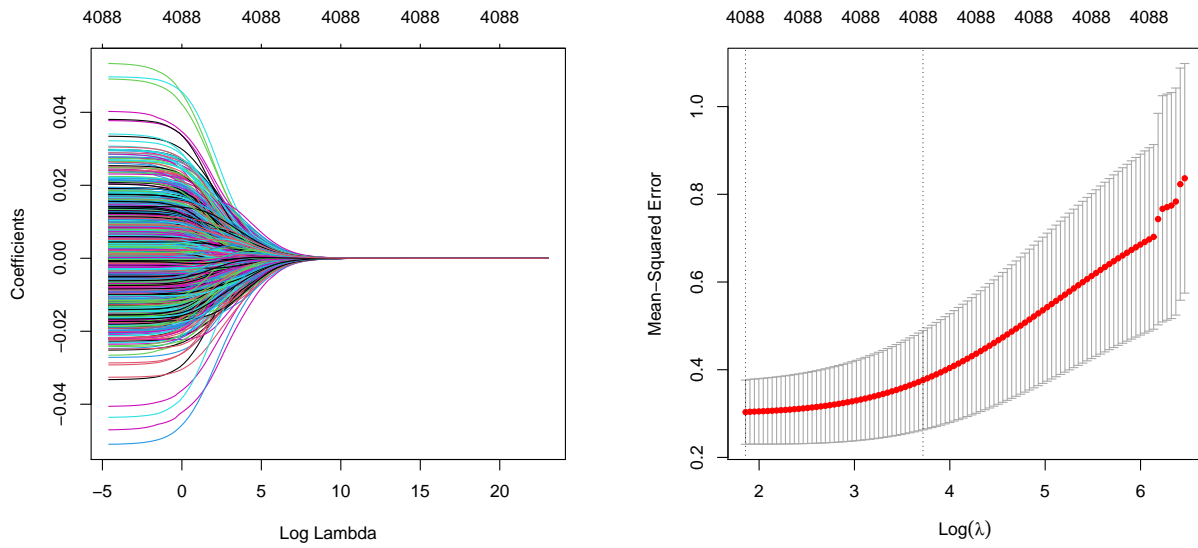
لگاریتم سطح ژن‌ها است.

برای تحلیل این داده‌ها ابتدا از روش مؤلفه اصلی استفاده می‌کنیم. نمودار لگاریتم بازو حاصل در شکل ۱۳ قابل مشاهده است. در ادامه از روش‌های کمترین توان‌های دوم تاوانیده استفاده می‌کنیم. برای تعیین پارامتر جریمه از روش اعتبارسنجی متقابل  $10^\circ$  -سطحی با در نظر گرفتن  $70\%$  داده‌ها به عنوان مجموعه آموزش و  $30\%$  به عنوان مجموعه آزمون استفاده می‌شود. نمودارهای مربوط به اعتبارسنجی و برآورد ضرایب به ازای پارامترهای تاوانیده مختلف با روش ستیغی و لاسو مشهود است. در نمودار اعتبارسنجی روش ستیغی به ازای مقدار  $67409536$  مقدار MSE را کمینه کرده و در روش لاسو مقدار آن به ازای  $19790790$  کمینه شده است. در قسمت فوقانی نمودارهای سمت راست شکل‌های ۱۴ و ۱۵ تعداد ضرایب غیر صفر موجود به ازای مقادیر مختلف  $\lambda$  نمایش داده شده است که واضح است که با افزایش میزان  $\lambda$  تعداد ضرایب غیر صفر روش لاسو کاهش می‌یابد درحالی که در روش ستیغی برآورد ضرایب بسیار به صفر نزدیک می‌شوند ولی دقیقاً معادل با صفر نخواهند بود. نمودار اعتبارسنجی برای روش شبکه ارتجاعی در شکل ۱۶ نمایش داده شده است. نمودارهای برآورد ضرایب به روش شبکه ارتجاعی در شکل ۱۷ رسم شده است. نمودارهای برآورد ضرایب به روش اسکد به ازای دو مقدار  $a = 3,37$  و بردار پارامتر دلخواه یکسان در شکل برای  $\lambda$  ۱۸ رسم شده است. نمودار اعتبارسنجی متقابل برای یافت مقدار بهینه پارامتر  $\lambda$  به ازای مقدار پیشنهادی فن و لی [۱۲] برای پارامتر  $a$  در شکل ۱۹ نشان داده شده است. در جدول ۳ نتایج روش‌های مورد استفاده گزارش شده است.

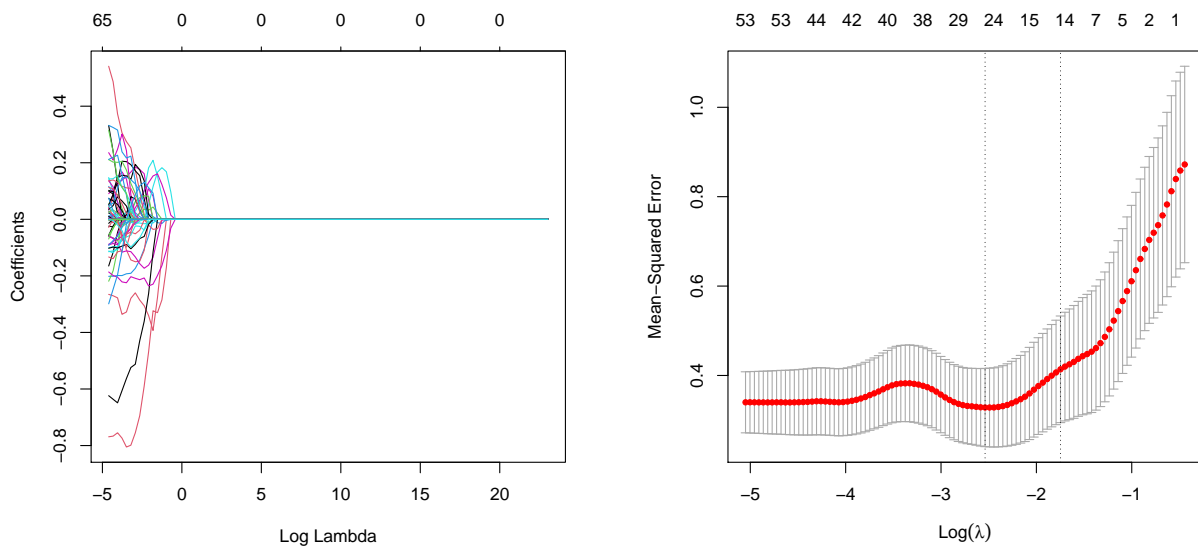
بعضی از محصولات غذایی اضافه می‌شود و به عنوان یک مکمل رژیم غذایی در دسترس است. این ویتامین یکی از اجزای اساسی دو کوآنزیم اصلی، فلاوین مونوکلوئید (FMN) و فلاوین آدنین دینوکلوئید (FAD) است. این کوآنزیم‌ها نقش عمده‌ای در تولید انرژی، عملکرد سلولی، رشد و نمو، و متابولیسم چربی‌ها، داروها و استروئیدها دارند. علاوه بر این، ریوفلاوین به حفظ سطح طبیعی هموسیستین، یک اسید آمینه در خون کمک می‌کند. بیشترین درصد ریوفلاوین پس از مصرف مواد گوشتی، سبزیجات، لبنیات و مقداری کمتر در مغزها و تخمه‌ها، حبوبات و سبوس غلات موجود توسط بدن دریافت می‌شود. بیشتر ریوفلاوین توسط روده کوچک جذب و مقدار کمی در کبد، قلب و کلیه‌ها ذخیره شده و مابقی از طریق ادرار از بدن دفع خواهد شد. وضعیت ریوفلاوین در افراد سالم به طور معمول اندازه‌گیری نمی‌شود. اما در صورت بروز مشکل و تشخیص کمبود (یا وفور) این ویتامین از طریق بررسی گلبول‌های قرمز خون و یا بررسی میزان ریوفلاوین دفع شده از طریق ادرار اندازه‌گیری خواهد شد. میزان میانگین ریوفلاوین مورد نیاز بدن روزانه بین  $1/1$  تا  $1/3$  میلی‌گرم برای بزرگسالان و برای زنان باردار یا شیرده  $1/6$  می‌باشد. شایع‌ترین علت کمبود ریوفلاوین در بدن رژیم غذایی نامناسب است. کمبود ریوفلاوین هم‌چنین می‌تواند در افرادی که دچار نقص در فعالیت‌های کبدی هستند ایجاد شود، چراکه مانع از استفاده مناسب از ویتامین‌ها می‌شود. این داده‌ها در بسته نرم‌افزاری "hdi" نرم‌افزار R موجود است. در این مجموعه داده متغیر پاسخ لگاریتم نرخ تولید ریوفلاوین است. این مجموعه داده دارای  $4088$  متغیر توضیحی بوده که هر کدام نشان‌دهنده



شکل ۱۳: لگاریتم نمودار بازو برای داده ریوفلاوین.



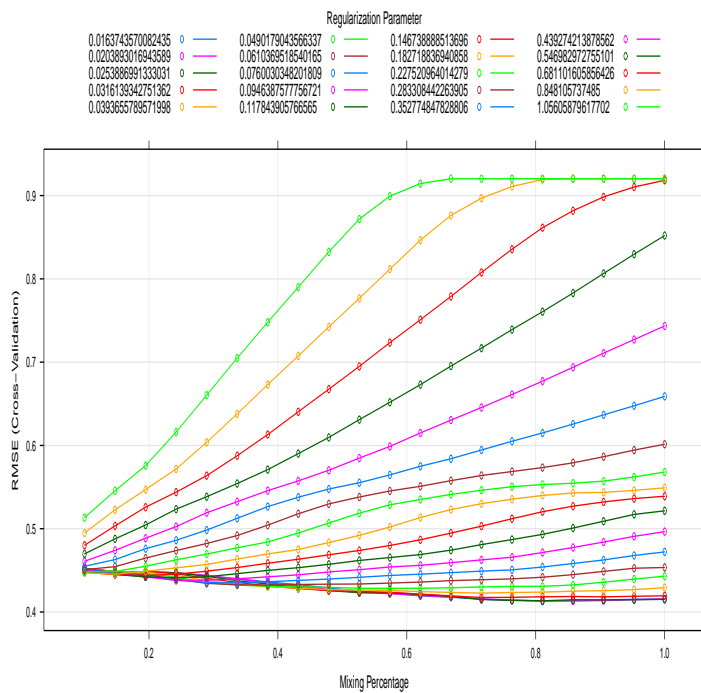
شکل ۱۴: نمودار اعتبارسنجی و برآورد ضرایب داده ریو فلاوین به ازای پارامترهای تاوان مختلف به روش ستیغی.



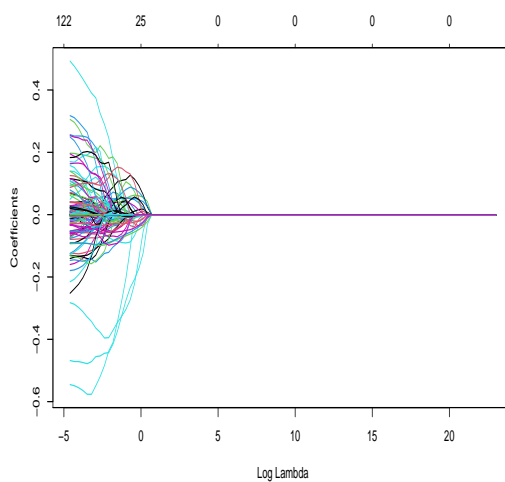
شکل ۱۵: نمودار اعتبارسنجی و برآورد ضرایب داده ریو فلاوین به ازای پارامترهای تاوان مختلف به روش لاسو.

جدول ۳: جدول نتایج روش‌های تاوانیده برای داده ریو فلاوین.

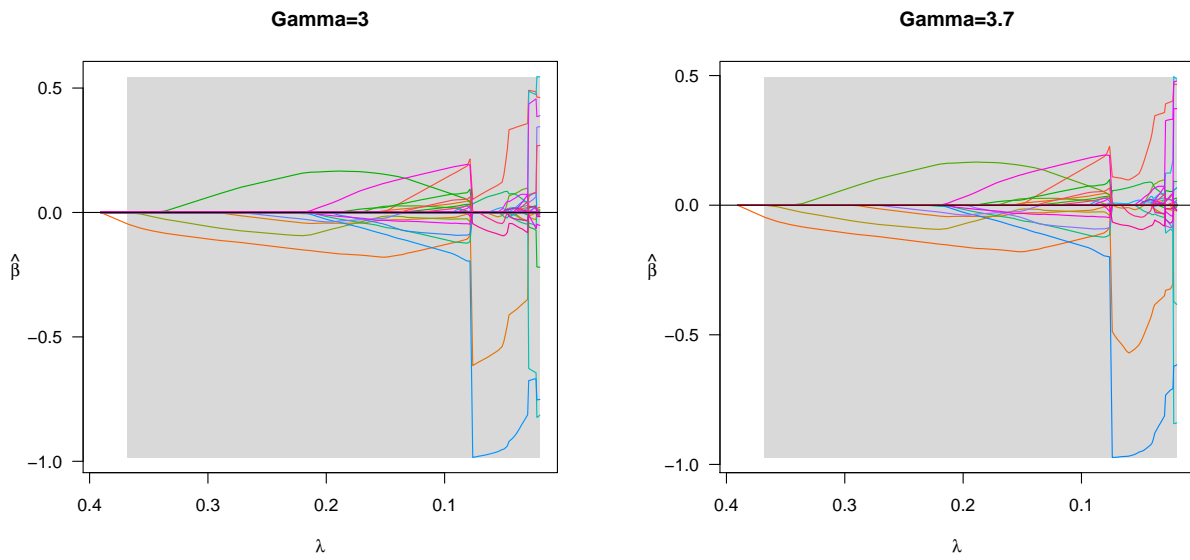
روش	مقدار بهینه پارامتر تاوان	تعداد ضرایب غیر صفر	مجموع مربعات خطا
ستیغی	۶۴۰۹۵۳۶	۴۰۸۸	۸۴۶۱۹۳۳
لاسو	۰٫۰۷۹۰۱۹	۲۸	۱۰۲٫۶۴۴
شبکه ارتجاعی	$\alpha = ۰٫۸۱۰۵۲, \lambda = ۰٫۰۲۰۳۸$	۴۳۷	۴۰۴۹٫۶۰۵
اسکد	$a = ۳٫۷, \lambda = ۰٫۰۵۷۱۵۱۱۳$	۱۰	۲٫۰۶۸۴۶۱



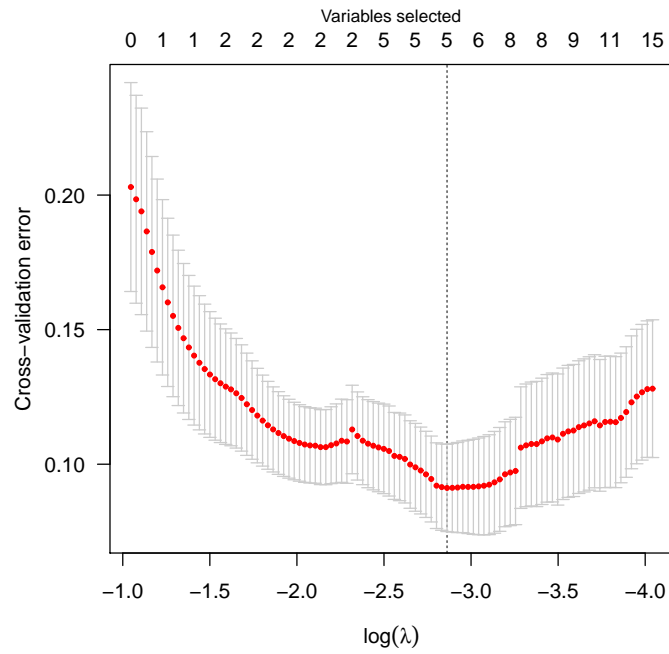
شکل ۱۶: نمودار اعتبارسنجی برای داده ریوفلاوین به روش شبکه ارتجاعی.



شکل ۱۷: نمودار ضرایب برای داده ریوفلاوین به روش شبکه ارتجاعی.



شکل ۱۸: نمودارهای روش اسکید برای داده ریوفلاوین به ازای  $\gamma = 3$  و  $\gamma = 3.7$ .



شکل ۱۹: نمودار اعتبارسنجی برای یافت پارامتر  $\lambda$  برای داده ریوفلاوین به روش اسکید.

## ۵ نتیجه گیری

برای انتخاب متغیر نیست (برخلاف روش لاسو). به همین دلیل هدف مقایسه روش‌های مطرح شده در این تحقیق نبوده و بیشتر ایجاد یک دید جامع در خوانندگان در مواجهه با این گونه داده‌ها است. بنابراین هدف اصلی تنها معرفی و به کارگیری هم‌زمان روش‌ها در داده‌های واقعی و شبیه‌سازی شده می‌باشد. اگرچه به کارگیری هم‌زمان چندین روش بسیار راهگشا خواهد بود و دید جامعی نسبت به انواع مشاهدات در اختیار محققان و کاربران قرار می‌دهد.

مشکل همخطی و یا بزرگ‌تر بودن تعداد متغیرهای توضیحی نسبت به مشاهدات، جواب‌های روش کمترین توان‌های دوم را به طرز قابل توجهی منحرف یا ناممکن می‌کنند. برای حل این مشکل در مسائل مدل‌سازی رگرسیون تاکنون روش‌های متعددی پیشنهاد و بررسی شده است که در این مقاله به معرفی مهم‌ترین آن‌ها پرداخته شد (از جمله سایر روش‌ها می‌توان به [۵] و [۲۱] اشاره کرد). در حالت کلی نمی‌توان گفت که کدام روش بر سایر روش‌ها برتری دارد، زیرا هر روش طبق شرایط داده‌های جمع‌آوری شده بهتر عمل می‌کند. مثلاً روش لاسو نمی‌تواند در مواردی که بین داده‌ها هم خطی وجود دارد خوب عمل کند ولی روش شبکه ارتجاعی در این شرایط بهتر عمل می‌کند. از طرف دیگر روش رگرسیون ستیغی اگرچه می‌تواند مشکل وجود همخطی بین متغیرهای توضیحی را حل کند، ولی روش مناسبی

## تقدیر و تشکر

نویسندگان مقاله ضمن تشکر از اعضای محترم هیئت تحریریه مجله، از پیشنهادها و نظرات ارزشمند داوران و ویراستار محترم مقاله که موجب ارتقاء سطح آن گردید کمال تشکر و قدردانی را دارند.

## مراجع

- [۱] آرشی، م.، صادقی، ح. و طباطبایی، م. (۱۳۹۸)، استنباط آماری، انتشارات دانشگاه صنعتی شاهرود، شاهرود.
- [۲] رستا، ر.، چینی پرداز، ر. و راسخی، ع.ا. (۱۳۸۹)، معرفی روش کمترین توان‌های دوم تاوان داده در انتخاب متغیرهای توضیحی، ششمین همایش ملی آمار، دانشگاه پیام‌نور، اهواز.
- [۳] جذن، س. و امینی، م. (۱۳۹۶)، برآورد استوار نسبت به مشاهده‌های دورافتاده در رگرسیون خطی در حضور هم‌خطی چندگانه، مجله اندیشه آماری، ۴۴، ۹۳-۱۱۰.
- [۴] نوری جلیانی، ک.، نوری، س.، محمد، ک.، نیکنام، م.ح.، محمودی، م.، آندونیان، ل. و اکبری، آ. (۱۳۹۰). آنالیز جنگل‌های تصادفی: یک روش آماری مدرن برای غربالگری در مطالعات با بعد بالا و کاربرد آن در یک مطالعه همبستگی ژنتیکی جمعیت-پایه، مجله دانشگاه علوم پزشکی خراسان شمالی، (ویژه نامه آمار زیستی و اپیدمیولوژی)، ۳، ۹۳-۱۰۱.
- [5] Akdeniz, F and Roozbeh, M. (2019). Generalized difference-based weighted mixed almost unbiased ridge estimator in partially linear models, *Statistical Papers*, **60(5)**, 1717-1739.
- [6] Bellman, R. (1961). *Adaptive control processes*. Princeton university press, London.
- [7] Bertsimas, D and Parys, B. V. (2020). Sparse high-dimensional regression: exact scalable algorithms and phase transitions, *Biostatistics*, **21(2)**, 219-235.
- [8] Bonddel, H.D. and Reich, B.J. (2008). Simultaneous regression shrinkage, variable selection and clustering of predictors with OSCAR, *journal of the interntional biometric society*. **64(1)**, 115-123.
- [9] Breiman L. (1995). Better subset regression using the nonnegative garrote, *Technometrics*. **37(4)**, 373-384.
- [10] Efron, B. and Hastie, T. (2017). *Computer age statistical inference*. Cambridge University Press, Cambridge.
- [11] Everitt, B. and Hothorn, T. (2011). *An introduction to applied multivariate analysis with R*. Springer, New York Dordrecht, Heidelberg, London.

- [12] Fan, J. and LI, R. (2001). Variable selection via nonconcave penalized likelihood and its oracle properties, *Journal of the American Statistical Association*. **96(456)**, 1348-1360.
- [13] Fan, J. (1997). Comments on wavelets in statistics: a review by A. Antoniadis,” *Journal of the Italian Statistical Society*. **6(2)**, 131-138.
- [14] Frank, E. and Friedman, J. (1993). A statistical view of some chemometrics regression tools, *Technometrics*. **35(3)**, 109-148.
- [15] Hastie, T. J. and Pregibon, (1992). *Generalized linear models*. Eberly College of Science, London.
- [16] Hoerl, A. E. (1962). Application of ridge analysis to regression problems, *Chemical Engineering Progress*. **58(1)**, 54-59.
- [17] Hoerl, A. E. and Kennard, R. W. (1970). Ridge regression: applications to nonorthogonal problems, *Technometrics*. **12(1)**, 69-82.
- [18] Jolliffe, I.T. (2002). *Principal component analysis*. Springer series in statistics, Aberdeen.
- [19] Li, B. and Yu, Q. (2009). Robust and sparse bridge regression, *Statistics and its interface*, **2(4)**, 481-491.
- [20] Pearson, K. (1901). On lines and planes of closest fit to systems of points in space, *Philosophical Magazine*. **2(1)**, 559-572.
- [21] Roozbeh, M. (2018). Optimal QR-based estimation in partially linear regression models with correlated errors using GCV criterion, *Computational Statistics & Data Analysis*. **117**, 45-61.
- [22] Roozbeh, M., Babaie-Kafaki, S. and Naeimi Sadigh, A. (2018). A heuristic approach to combat multicollinearity in least trimmed squares regression analysis, *Mathematical modelling*. **57(2)**, 105-120.
- [23] Tibshirani, R. (1996). Regression shrinkage and selection via the Lasso, *Journal of the Royal Statistical Society. Series B (Methodological)*. **58(1)**, 267-288.
- [24] Walker, D. A. and Smith, T. J. (2020). Logistic regression under sparse data conditions, *Journal of Modern Applied Statistical Methods*, **18(2)**, 33-72.
- [25] Wasserman, L. (2006). *All of nonparametric statistics*. Springer Science and Business Media, New York.
- [26] Watkins, D.S. (2002) *Fundamentals of matrix coputations*. John Wiley and Sons, New York.
- [27] Zou, H. and Hastie, T. (2005). Regularization and variable selection via the elastic net, *Journal of the Royal Statistical Society, Series B*. **67(2)**, 301-320.